

MEMORIAS DE LA XXVII ESCUELA DE VERANO EN FÍSICA
JUNIO 10 - 21, 2019

Eds: José Récamier

Rocío Jáuregui

Instituto de Física, UNAM, 10 - 14 de junio 2019
Instituto de Ciencias Físicas, UNAM, 17 - 21 de junio de 2019

Contenido

Contribuciones

Introducción

Agradecimientos

Profesores participantes

Alumnos participantes

- **Contribuciones**

Antonio M. Juárez Reyes	1
<i>Aplicaciones recientes de espectroscopia en problemas relevantes a la sociedad</i>	
Sebastien Fromenteau	12
<i>Numerical and Statistics methods for Cosmology</i>	
Guillermo Hinojosa Aguirre	24
<i>La física fundamental y aplicaciones de los iones negativos</i>	
François Leyvraz	34
<i>Entrelazamiento, no-localidad y otras particularidades de la mecánica cuántica.</i>	
Frédéric S. Masset	45
<i>Pplanetary migration in gaseous protoplanetary disks</i>	
I Ramos-Prieto, A Paredes, H. Moya-Cessa y J. Récamier	64
<i>Evolución aproximada para un sistema compuesto por dos Hamiltonianos de Jaynes-Cummings acoplados</i>	
Thomas Stegmann	77
<i>A brief introduction to the NEGF method for electron transport at the nanoscale</i>	
J Alberto Vázquez, R. Medel Esquivel, I Gómez Vargas	89
<i>Cosmología observacional con Redes Neuronales Artificiales</i>	
María Barrera, Carlos Villarreal	95
<i>Redes reguladoras complejas en la Medicina y la Biología</i>	

- **Introducción**

La XXVII Escuela de Verano en Física fue organizada con apoyo del Posgrado en Ciencias Físicas por el Instituto de Física y el Instituto de Ciencias Físicas de la Universidad Nacional Autónoma de México. Se llevó a cabo en las instalaciones del Instituto de Física en Ciudad Universitaria, del 10 al 14 de junio de 2019 y en las instalaciones del Instituto de Ciencias Físicas, en Cuernavaca, Morelos, del 17 al 21 de junio de 2019.

En esta Escuela se impartieron 18 cursos y 11 conferencias. Los cursos y conferencias cubrieron un amplio espectro con temas como óptica cuántica, sistemas complejos, colisiones atómicas y moleculares, mecánica cuántica, química cuántica, movimientos colectivos, cosmología y biofísica, entre otros.

Rocío Jáuregui, IF UNAM
José Récamier, ICF UNAM
Febrero, 2021

- **Agradecimientos**

Agradecemos los apoyos recibidos para la realización de esta escuela a la Universidad Nacional Autónoma de México a través de la Coordinación de la Investigación Científica, del Instituto de Física y del Instituto de Ciencias Físicas.

- **Profesores participantes**

- Maximino Aldana González, ICF UNAM, *Modelos matemáticos de la simbiosis entre bacterias y humanos*
- Iván Ortega Blake, ICF UNAM, *Biofísica molecular*.
- Horacio Martínez Valencia, ICF UNAM, *Plasmas y aplicaciones*.
- José Récamier Angelini, ICF UNAM, *Optomecánica*
- Antonio M Juárez Reyes, ICF UNAM, *Aplicaciones recientes de espectroscopia en problemas relevantes a la sociedad*.
- Kurt Bernardo Wolf Bogner, ICF UNAM, *Mecánica cuántica finita y sus aplicaciones en análisis de señales*
- Juan Carlos Degollado, ICF UNAM, *Relatividad numérica*
- José Alberto Vázquez, ICF UNAM, *Cosmología observacional con Redes Neuronales Artificiales*
- Guillermo Hinojosa Aguirre, ICF UNAM, *La física fundamental y aplicaciones de los iones negativos*
- Thomas Stegmann, ICF UNAM, *Transporte electrónico en nano-sistemas*
- Víctor Contreras, ICF UNAM, *Espectroscopia atómica de líquidos acústicamente levitados*
- Juan Carlos Hidalgo Cuéllar, ICF UNAM, *Estadística Bayesiana y su aplicación a la cosmología*
- W Luis Mochán, ICF UNAM, *Metamateriales*
- Carlos Villarreal Luján, IF UNAM, *Redes complejas, salud y enfermedad*
- Luis Benet Fernández, ICF UNAM, *Integración de EDOs precisa e integración validada*.
- Hernán Larralde Ridaura, ICF UNAM, *Caminatas aleatorias*.
- Sebastien Fromenteau, ICF UNAM, *Cosmología observacional: de la teoría a la observación gracias a los campos aleatorios Gaussianos*
- François Leyvraz Waltz, ICF UNAM, *Entrelazamiento, no-localidad y otras extrañezas de la mecánica cuántica*

- Agustín González Flores, ICF UNAM, *Algunos aspectos de la cristalización coloidal en 2D.*
- Rocío Jáuregui Renaud, IF UNAM, *Correlaciones cuánticas*
- Carlos Pineda, IF UNAM, *Sistemas cuánticos abiertos*
- Freddy Jackson, IF UNAM, *Gases fríos y ultrafríos: Plataformas a la simulación cuántica*
- Alejandro Pérez Riascos, IF UNAM, *Procesos dinámicos en sistemas complejos modelados por redes.*
- Luis Acosta, IF UNAM, *Física nuclear en astrofísica.*
- Oliver Paz, Andrés Botello y Penélope Rodríguez, IF UNAM, *Simulación atomística de materiales*
- Aurora Courtoy, IF UNAM, *Protones, neutrones, piones, etcétera.*
- Ricardo Méndez Fragoso, FC UNAM, *Soluciones analíticas y numéricas de la ecuación no-lineal de Schrödinger.*
- Rubén Alfaro, IF UNAM, *Observando el Universo a muy altas energías.*
- Erick Vázquez, IF UNAM, *Física de astropartículas: detectando neutrinos y buscando la materia oscura del Universo*

- **Alumnos participantes**

- José David Romo López (FC-UNAM) davidrl96@ciencias.unam.mx
- Torres Arvizu Francisco Ricardo (FC-UNAM) toaf951023@ciencias.unam.mx
- Manuel Alejandro Alderete (FC-UNAM) m_alderete_lezama@ciencias.unam.mx
- Aarón Ali Lozano Rocha (FC-UNAM) aliensnero@gmail.com
- Fausto Felipe Morales González (FC-UNAM) morales@ciencias.unam.mx
- Andrea Pizarro Medina (FC-UNAM) a.pizarro@ciencias.unam.mx
- Alejandro Pérez Fernández(FC-UNAM) alex.perezfer@ciencias.unam.mx
- Cristian Carvajal Bohorquez (Colombia) cristian_9807@hotmail.com
- Juan Antonio Saldaña Herrera(FC-UANL) jasaldanah@gmail.com
- María Magdalena Castro Sam (BUAP) ma.castro.sam@gmail.com
- María Guadalupe Morales Trejo (UAQ) lu.moralestrejo@hotmail.com
- Juan Felipe Pulgarin Mosquera (Colombia) jfelipe.pulgarin@udea.edu.co
- Lizeth Daniela Jaimes González(Colombia) danielajaimesg15@gmail.com
- William Andrés Jaimes Espíndola (Colombia) william.jaimes01@gmail.com
- Juan José Segura Flóres (Colombia) juan.segura01@uptc.edu.co
- Mariana Jaramillo Acero (Colombia) mariana.jaramillo2@udea.edu.co
- Valeria Itzel Arteaga Muñoz (UAZ) valitzelar@hotmail.com
- José Gustavo Bravo Flores (UNISON) jgustavobflores@gmail.com
- David Ramos Salamanca (Colombia) david.ramos.salamanca@outlook.com
- Lina María Montoya Zuluaga (Colombia) linam.montoya@udea.edu.com
- Zamir David Beleño Rodríguez (Colombia) zbeleno@mail.uniatlantico.edu.co
- Alberto Pedraza Pedraza (BUAP)washburnuser@gmail.com
- Francisco Valentín Valerio López (BUAP) trainwerck1979@gmail.com
- Luis Angel Méndez López (FC UNAM)
- Aniel de Jesús Villegas (FC UNAM)

Aplicaciones recientes de espectroscopia en problemas relevantes a la sociedad

Dr. Antonio M. Juárez

Laboratorio de Fotodinámica

Instituto de Ciencias Físicas, UNAM

La ciencia, tanto fundamental como aplicada, es el pilar en el que se apoyan las sociedades avanzadas y con buen desarrollo económico a nivel mundial. La relación de beneficio directo que tiene a la ciencia con el bienestar de la sociedad que la cultiva es incluso cuantificable empleando el índice de economía del conocimiento. En México, a pesar de ser un país grande y con un gran potencial económico la ciencia, ya sea esta aplicada o fundamental representa un porcentaje muy pequeño, menor al 0.4% en relación con el producto interno bruto del país. Lo anterior es, al mismo tiempo causa y efecto del hecho que nuestro país dedica su labor económica a servicios, producción de materia prima sin valor agregado alto y empleo en sectores de manufactura. La solución a este tipo de dificultades es compleja e implica varios factores que van desde la poca cultura científica del país, la baja inversión tanto privada como pública para labores de investigación y el bajo ingreso de alumnos a carreras de base científica. Desde la academia, los que dedicamos nuestros esfuerzos a desarrollar labor científica, es posible contribuir en el último aspecto que se mencionó: Fomentar vocaciones científicas. Una manera concreta de realizar esto es difundiendo de la manera más explícita posible el vínculo profundo y directo que tiene la ciencia en la generación de riqueza y bienestar, además de ser un valor cultural por si misma.

Esta contribución tiene como propósito mencionar diversas aplicaciones que tiene la física molecular y la espectroscopia en diversos aspectos prácticos de la vida diaria y de importancia estratégica para la sociedad mexicana. El propósito de presentar esta lista extensa es el de motivar en los estudiantes el interés por cultivar estas áreas a la vez que apreciar el impacto directo que tienen los desarrollos científicos de alto nivel en aspectos tan centrales para la sociedad como lo son las telecomunicaciones, la salud, la seguridad y defensa nacional y la seguridad alimentaria. No es exagerado afirmar que, prácticamente en todos los aspectos centrales de la sociedad, hay una contribución importante de la física. Tanto las distintas técnicas de espectroscopia como las aplicaciones de estas son tan bastas que es necesario acotar el alcance de esta contribución. Concretamente, se mencionarán en este documento las aplicaciones de la espectroscopia de absorción o de fluorescencia en los rangos del ultravioleta, el visible y el infrarrojo cercano. Existen muchas otras técnicas como la espectroscopia de masas, espectrometría de rayos X, espectroscopia Raman, espectroscopia FTIR y varias más.

Espectroscopia VIS-UV y NIR, aspectos básicos.

En primer lugar, es conveniente describir los aspectos fundamentales de las técnicas de espectroscopia en el visible, ultravioleta e infrarrojo que se abrevian como Vis-UV y NIR, respectivamente y por sus siglas en inglés. En primer lugar, la espectroscopia ultravioleta-

visible se realiza en el rango comprendido entre los 190 y 800 nanómetros de longitud de onda. Este tipo de espectroscopía se lleva a cabo empleando una fuente de luz de espectro amplio, típicamente lámparas de tungsteno o excímeros y un espectrofotómetro. La fuente de luz se hace incidir sobre una muestra y el espectrofotómetro analiza el espectro de reflexión. La luz que se hace incidir en la muestra se absorbe de manera selectiva, dependiendo de los grupos funcionales o de la estructura molecular, única, de la muestra bajo estudio. Esta luz absorbida promueve, a nivel interno del material la excitación de electrones a bandas de energía permitidas, o bien de transiciones vibracionales o rotacionales, en el caso de la espectroscopía NIR. comprende entre 190 y 800 nm.

Uno de los instrumentos más versátiles y empleados para realizar tanto la espectroscopía NIR como VIS-UV es el espectrofotómetro en la geometría de *Zerny-Turner*. El monocromador *Zerny-Turner* es un dispositivo óptico que sirve para medir la composición de la luz según su distribución de longitudes de onda diseñado por *M. Czerny* y *A.F. Turner* en 1930. Este dispositivo consta de 6 elementos ópticos: (A) un iris (comúnmente una fibra óptica) por el cual entra un haz de luz policromática que apunta a una (B) ranura. La ranura regula la intensidad de luz que entra en el dispositivo se determina por la intensidad de la fuente y las dimensiones de la ranura (base X altura). La ranura se coloca en el punto focal efectivo de un espejo cóncavo (el colimador (C)) que colima el haz de luz de forma tal los rayos incidan paralelos sobre la rejilla de difracción (D). La rejilla separa el haz de luz en sus distintas componentes espectrales las, enviando cada componente espectral con en un ángulo distinto. Estos haces de luz se envían a otro espejo cóncavo (el de enfoque (E)) que enfoca el haz sobre una salida (F), que puede ser una ranura, una pantalla o un dispositivo analógico digital (por ejemplo, un CCD lineal). En la Figura 1 se pueden ver las distintas componentes ópticas del sistema de manera esquemática.

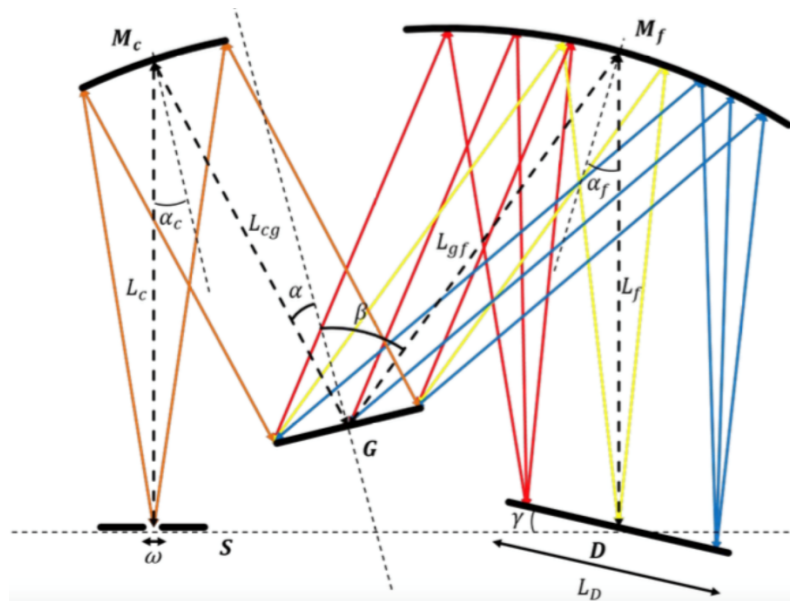


Figura 1. Representación esquemática de un espectrofotómetro *Zerny-Turner*, empleado en espectroscopía VIS-UV y NIR.

Una vez presentados los aspectos básicos de la técnica de espectroscopía VIS-UV se presentan a continuación aplicaciones directas de esta técnica en diversas aplicaciones.

Aplicaciones en seguridad alimentaria.

La inocuidad de los alimentos, que se define como la ausencia de sustancias o microorganismos en los alimentos es uno de los aspectos más relevantes, actualmente, en relación a la salud pública y la sustentabilidad de las sociedades modernas. Con el fin de asegurar esta inocuidad, la industria alimentaria emplea diversas técnicas para asegurar que los alimentos que llegan a los consumidores no presenten pesticidas, sustancias tóxicas o microorganismos. La espectroscopía VIS UV presenta, actualmente, desarrollos muy prácticos y de costo relativamente bajo para llevar a cabo este tipo de análisis. Aunque se han desarrollado numerosas estrategias para detectar, aislar e identificar posibles amenazas en los alimentos, es de fundamental importancia contar con técnicas que mejoren la velocidad, la sensibilidad y la selectividad de la detección de compuestos tóxicos. Esto se debe lograr al mismo tiempo que las técnicas que se empleen sean simples, portátiles y de bajo costo. Una de las debilidades de la técnica VIS-UV es la sensibilidad limitada que tienen los espectrofotómetros para determinar la presencia de compuestos que se encuentren diluidos en proporciones muy bajas.

Las normas internacionales, por otro lado, piden, justamente, que la detección de pesticidas o compuestos nocivos se lleve a cabo a niveles de detección de partes por millón o incluso menores. Para subsanar esta dificultad de detectar trazas de compuesto tóxicos, se han desarrollado materiales que actúan como amplificadores de señales. Uno de estos son los puntos cuánticos. Estos compuestos se componen de nanocristales semiconductores fluorescentes brillantes. La fluorescencia de estos compuestos puede ajustarse para atenuar o amplificar señales específicas de compuestos. Existen actualmente diversas investigaciones sobre la utilización de puntos cuánticos en ensayos para la detección de analitos como patógenos, pesticidas, antibióticos e incluso de organismos genéticamente modificados (*Xiaon Lu*, 2017). También se incluye una breve introducción sobre las propiedades y la bioconjugación de puntos cuánticos. Numerosos estudios han demostrado el potencial de los puntos cuánticos para mejorar las cifras analíticas de mérito en ensayos de calidad y seguridad alimentaria; sin embargo, se necesita investigación estratégica para desarrollar ensayos con puntos cuánticos que tengan la mayor oportunidad de impactar las prácticas de seguridad alimentaria en la industria y la sociedad. El campo es extensísimo, por lo que se sugiere al lector interesado consultar las referencias relevantes (*Xingbo Xi et al*, 2019). Otra área de particular relevancia en la seguridad alimentaria es el uso de espectrofotometría en el infrarrojo cercano. Esta técnica, que explora la luz dispersada por alimentos en el rango de 800 hasta 2000 nanómetros tiene una gran cantidad de aplicaciones (*Jia-Huan Qu*, 2015). Aunque, evidentemente, existen procedimientos de control de calidad y seguridad utilizando técnicas tradicionales, como métodos instrumentales y fisicoquímicos y procedimientos microbiológicos, estos procedimientos son lentos y caros, por lo que una alternativa

espectroscópica puede ser de gran utilidad. El rango de espectroscopía NIR proporciona información directa de los modos de interacción que proporcionan complejos químicos y físicos información relacionada con el comportamiento de vibración de molecular de enlaces tales como C-H, O-H, y N-H. Las aplicaciones concretas de esta técnica son muy variadas y mencionaremos solo algunas selectas. Usando espectroscopía NIR es posible determinar el nivel de frescura de carnes, así como la potencial adulteración de miel, vinos, licores o jugos. De particular interés para nuestro país es la exportación de el aguacate. El valor de este producto en el mercado norteamericano es muy elevado y, anualmente, se exportan cantidades altas de este producto que representan un mercado de casi 3 mil millones de dólares anuales. Una de las propiedades más importantes que determinan la calidad de los aguacates de exportación es el contenido de materia seca. La espectroscopía NIR se puede aplicar en determinar esta importante propiedad, con el reto agregado de que debe determinarla en 10 milisegundos, que es el tiempo que transita la fruta frente a los espectrómetros, en un proceso industrial de empaque (*C.J Clark, 2003*).

Aplicaciones en monitoreo en la industria microelectrónica.

La detección y cuantificación de especies gaseosas en concentraciones muy bajas es de vital interés en varias áreas de la ciencia y la industria. En el entendimiento de los procesos de la química de la atmósfera o para mantener altos rendimientos en procesos industriales con materiales ambientalmente sensibles como los que encontramos en microelectrónica, o incluso en la seguridad del aire que respiramos, necesitamos de instrumentos analíticos altamente sensibles, confiables, y selectivos, para contar con información en tiempo real de qué moléculas gaseosas están presentes en cierto entorno y en cuánta cantidad.

Desde el punto de vista industrial, existe la necesidad de medir ciertos gases a niveles muy bajos de dilución, ya sea para mantener altos rendimientos del proceso o resultados, y para que la calidad de los procesos sea consistente, así como para la seguridad del propio personal de trabajo (*P. R. Griffiths, 2007*), (*Siefering, 1993*), (*Lebens et al, 1996*), (*Saga and Hattori, 1997*) (*Kawai, 1994*). Por ejemplo, en la fabricación de semiconductores se debe eliminar el amoniaco (NH₃) del aire ambiente ya que puede afectar negativamente los rendimientos en los procesos de fabricación de obleas incluso en concentraciones por debajo de una parte por millón (ppm).

Alternativamente, monitorear en tiempo real la acumulación de gases peligrosos para la vida humana como el monóxido de carbono (CO), óxidos de nitrógeno (NO_x), cloruro de hidrógeno (HCl) y fluoruro de hidrógeno (HF), son esenciales para la seguridad de los empleados. Por ejemplo, el HF es un químico altamente tóxico y corrosivo que también se encuentra comúnmente en muchas otras industrias menos prístinas, como la fundición de aluminio y la fabricación de artículos de vidrio, donde debe ser eliminado del aire ambiente por razones de seguridad en la salud humana. Cuando las dimensiones críticas del dispositivo semiconductor caen por debajo del rango de un cuarto de micrómetro, el concepto de ambiente libre de partículas toma especial relevancia (*Benjamin, 1993*). Por lo que a la transformación en la configuración de los cuartos limpios, se le sumó la

introducción de la unidad de filtro de ventilador (FFU, del inglés *fan-filter unit*) que consiste en un equipo motorizado para filtrar las partículas dañinas en el aire, y la implementación de la interfaz mecánica estándar (SMIF, del inglés *standard mechanical interface*) que es una pequeña estructura protectora para aislar las obleas de la contaminación creando un ambiente miniatura con un flujo de aire controlado. A pesar de estos esfuerzos precautorios y, a veces costosos, el deterioro ocasional e irregular de las características físicas de las obleas se conservaba por razones desconocidas, hasta que los defectos se detectaron y se vincularon con los entornos de las instalaciones. Este tipo de micro-contaminación se conoce como contaminación molecular en el aire (AMC, del inglés *airborne molecular contaminants*), que es un término bastante genérico porque los contaminantes pueden estar en forma de gas, vapor o incluso aerosoles con naturalezas químicas muy diferentes. Por lo tanto, aunque el control de estos AMCs ha sido reconocido como un requisito de diseño esencial para todas las nuevas instalaciones de fabricación de semiconductores, la complejidad del problema, como la variación en la fuente y la concentración ambiental, así como los impactos sobre el proceso de fabricación, en gran medida obstaculizan el desarrollo de una estrategia de control efectiva o estandarizada.

Por lo tanto, no solo las partículas serán un factor difícil de controlar en un cuarto limpio sino también las AMCs, cuya tasa de llegada a la superficie de la oblea es varios órdenes de magnitud mayor que la de las partículas (*Amy H.K*, 1994). En la actualidad, contamos con datos sobre la influencia de la contaminación orgánica molecular en los procesos de semiconductores, y a partir de estos datos se desarrollaron las recomendaciones para los niveles críticos de contaminación (*Budde et al*, 1995). Se descubrió que la presencia de compuestos orgánicos en las obleas antes de la formación de la compuerta de óxido (conocido comúnmente como *oxide-gate*) reduce significativamente la calidad del óxido (*Iwamoto*, 1997). También, se observó un dopaje involuntario al desgasificar fósforo que contenía retardantes de llama, y se modificó el tiempo de incubación en la deposición de vapores químicos a baja presión de nitruro de silicio sobre sustratos de silicio. Sin embargo, los gases inorgánicos también pueden generar precipitados dañinos, causando daños irreversibles a las obleas procesadas. El organismo internacional para guiar a la industria de los semiconductores conocido como ITRS (*International Technology Roadmap for Semiconductors*) en el 2005 destacó la necesidad de mayores controles para las AMCs, por lo que proveedores de equipos y los fabricantes de circuitos integrados han financiado gastos enormes para la identificación y control de estos contaminantes que afectan el rendimiento del circuito integrado.

Una excelente opción para la detección y cuantificación de AMCs, es la Espectroscopia de Absorción de Banda Ancha Estimulada en Cavidades (BBCEAS, por sus siglas en inglés), la cual cuenta con una alta sensibilidad y resolución espectral debido al uso de cavidades ópticamente estables, y al utilizar un LED como fuente de luz combina la flexibilidad de una amplia ventana espectral que permite la detección de múltiples especies, con compactes en tamaño y bajo costo.

Aplicaciones en monitoreo de calidad de aire.

El aire que respiramos es uno de los aspectos más importantes de nuestra salud y existencia, así como la del resto de los seres vivos que nos sostenemos de éste. Debido a esta importancia, existen esfuerzos continuos para establecer reglamentación y controles para reducir la emisión de contaminantes. Un aspecto muy importante para implementar estas reglamentaciones consiste en el monitoreo de los gases, su cuantificación y la emisión de alertas, de ser el caso. Existen numerosas técnicas que permiten realizar este tipo de estudios en tiempo real, tales como el monitoreo remoto desde satélites, las imágenes hiperspectrales, entre otras. De éstas, sobresalen las técnicas espectroscópicas por su facilidad de implementación. Desafortunadamente, para determinar cuantitativamente la presencia de contaminantes, se requieren de técnicas muy sensibles, que permitan cuantificar la presencia de contaminantes en partes por mil millón volumétricas. En este sentido, existe una necesidad alta de equipos que determinen la presencia de óxidos de azufre, óxidos de nitrógeno, el ozono troposférico (a baja altura, diferente del estratosférico, que es muy beneficioso), el monóxido de carbono y los compuestos orgánicos volátiles. -

Recientemente, en el Instituto de Ciencias Físicas desarrollamos una técnica muy sensible para la detección de Dióxido de Nitrógeno (NO₂) que es uno de los contaminantes más dañinos en la atmósfera y nocivo para la salud humana.

La técnica que hemos desarrollado permite detectar compuestos con una alta sensibilidad para detectar especies en niveles de partes por millón (ppm) a partes por billón (ppb) o trillón (ppt) en el ambiente. También, esta técnica provee mediciones selectivas que no estén influenciadas por otras especies presentes en la muestra gaseosa a analizar. Esto es particularmente importante debido a la compleja y rápidamente variada mezcla de especies de trazas de gases. En el resto de esta sección se presentarán los principios de esta técnica, que actualmente está disponible en el Instituto de Ciencias Físicas. Se hará un énfasis especial en las aplicaciones de ésta con el fin de motivar en los lectores la curiosidad de su uso. A continuación, se hace una descripción básica de la técnica y sus antecedentes. Posterior a esta descripción se dará una glosa breve de las aplicaciones de esta técnica, desarrollada en el ICF-UNAM.

*En algunos trabajos se ha probado que en una cavidad óptica iluminada por un láser de banda angosta (5 – 100Hz), operando de forma continua, CW (del inglés, *Continuos Wave*), la intensidad de la luz transmitida a través de la cavidad en un estado estacionario es directamente proporcional al tiempo *ring-down* de la cavidad. Igualmente, para una cavidad iluminada con una fuente de luz CW de banda ancha, cada componente de longitud de onda alcanza su propio estado estacionario de intensidad, dependiendo de la intensidad de la fuente de luz y de los procesos de pérdida que afectan a los fotones en esa longitud de onda dentro la cavidad (donde lo último está claramente relacionado con el tiempo *ring-down*). En la práctica BBCEAS integra un espectro de la luz transmitida por la cavidad, dispersando esta luz en longitudes de onda mediante un espectrómetro y registrándola con un detector multi-elemento, por ejemplo, una cámara CCD o un arreglo de diodos lineales. El diagrama*

esquemático del montaje experimental de BBCEAS es presentado en la figura 2.4.

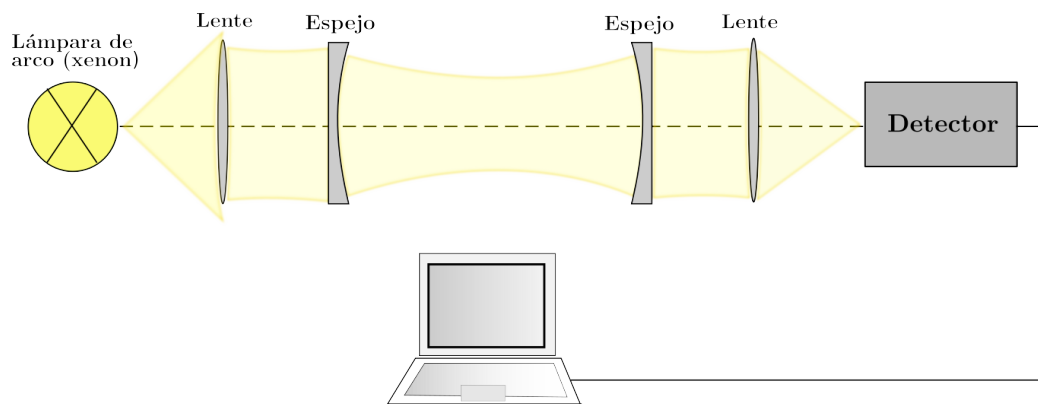


Figura 2.4: Esquema experimental de BBCEAS usando una lámpara de arco como fuente de luz.

iniciar o propagar la química en la atmósfera y, por lo tanto, son de gran interés en el estudio de la ciencia atmosférica. Muchas de las moléculas atmosféricamente importantes, previamente detectadas por la espectroscopía óptica de absorción diferencial, DOAS (del inglés, *Differential optical absorption spectroscopy*) en el visible y UV-cercano, ahora también pueden ser detectadas usando BBCEAS. La sensibilidad de ésta técnica es comparable con DOAS ya que ambas se basan en la medición de una señal de absorción de especies específicas en presencia de otros absorbentes y de aerosoles atmosféricos. La técnica BBCEAS tiene además la ventaja de poder manipular la muestra antes de entrar a la cavidad ring-down, permitiendo así detectar especies específicas que por sí mismas no absorben en longitudes de onda en el visible o en el ultravioleta cercano. En el 2004, *Ball S. M. et al.* presentaron una variante de BBCEAS en la región visible del espectro electromagnético usando LEDs de color rojo y verde y midiendo el espectro del oxígeno molecular y el vapor de agua. Adicionalmente, midieron los espectros de absorción de tres importantes absorbentes en la atmósfera: NO₃, NO₂, e I₂. En el 2006, *Dean S. Venables et al.* describieron la aplicación de BBCEAS en la detección *in situ* de trazas de gases y radicales atmosféricos (NO₃, NO₂, O₃, H₂O) en una cámara de simulación atmosférica bajo condiciones atmosféricas reales, alcanzando una sensibilidad de 4pptv para NO₃ en un tiempo de adquisición de 1 minuto. En el 2008 *Gherman T. et al.* reportaron la primera aplicación de BBCEAS en el ultravioleta-cercano para mediciones simultáneas de HONO y NO₂, logrando una sensibilidad de ~ 4ppbv para HONO y de ~ 14 ppbv para NO₂ en un tiempo de adquisición de 20 segundos. *Triki M. et al.* presentaron un arreglo experimental basado en un LED a 643 nm el cual es de interés en la detección simultánea de NO₃ y NO₂, con un límite de detección en el rango ppbv en un tiempo promedio de 2 minutos, comparable con la mejor versión de los dispositivos de quimioluminiscencia usados en el

monitoreo de contaminantes atmosféricos. En el 2013 *Liuyi Ling et al.* describieron la aplicación de esta técnica en mediciones *in situ* de NO₂ atmosférico usando un LED azul logrando una sensibilidad de 1 a 35 ppbv. Estudios recientes (*Washenfelder et al.*, 2016) en BBCEAS describen la aplicación de esta técnica en mediciones simultáneas de formaldehído (CH₂O) y NO₂ en la región ultravioleta de 315–350 nm, reportando un límite de detección de 300 pptv para CH₂O en un tiempo de adquisición de 1 minuto.

La técnica BBCEAS se emplea casi exclusivamente para la detección de gases que juegan un papel importante en la química de la atmósfera, con unas pocas excepciones en trabajos que reportan el uso de la técnica en mediciones de gases generados en procesos industriales, como es el caso del 1,3-butadieno (C₄H₆), un contaminante gaseoso peligroso producido a escala industrial para la generación de cauchos sintéticos y plásticos. *Denzer et al.* desarrollaron un instrumento BB-CEAS usando un diodo emisor de luz superluminescente en el IR-cercano para la detección de este contaminante, logrando una detección mínima del coeficiente de absorción de $6,1 \times 10^{-8} \text{cm}^{-1}$ en un tiempo de integración de pocos minutos. Otro contaminante presente en ambientes industriales es el 1,4-Dioxano (DX), el cual se usa ampliamente como solvente industrial en productos farmacéuticos y pinturas cosméticas. También se utiliza como agente humectante y dispersante en las industrias textiles y de tintes.

Cabe señalar que las configuraciones anteriores dependen de espejos con control micrométrico, los cuales se requieren para proporcionar la alineación de la cavidad. Esto a su vez, los hace propensos a la desalineación y reduce su portabilidad.

En este trabajo, se presenta un nuevo diseño de espectrómetro para la detección y cuantificación de trazas gaseosas en el rango de 610-670 nm, utilizando un LED como fuente de luz centrada a 634 nm. Este prototipo tiene un diseño monolítico (i.e. una sola pieza), en contraste con los diseños anteriores, lo que permite ensamblar el LED, las lentes, los espejos y el espectrómetro de baja resolución en una sola pieza de aluminio. La principal ventaja de este diseño es que no se requiere alinear los espejos antes de cada medida, lo que hace que el sistema sea práctico de ensamblar y usar, térmicamente estable y fácil de transportar. Además, se realizaron pruebas a diferentes temperaturas que muestran que la cavidad es robusta frente a cambios de temperatura de $\pm 10^\circ\text{C}$ alrededor de la temperatura ambiente de 26°C . Esto permite la medición de espectros en el campo, y no solamente en un entorno ambientalmente controlado.

Conclusiones.

La espectroscopia es un área de gran utilidad práctica en diversos campos aplicados. Este resumen tiene el propósito de dar algunos detalles de un número limitado de aplicaciones. Los alumnos interesados en profundizar en el tema son bienvenidos en el laboratorio de fotodinámica del Instituto de Ciencias Físicas UNAM o contactar directamente al autor de este artículo en el correo amjuarez@icf.unam.mx.

H. K. Amy J. Muller, Linda A. Psota-Kelty and J. Sinclair, "Volatile cleanroom contaminants: sources and detection," *Solid State Technology*, pp. 61–72, 1994

Y. Benjamin, Y. Liu, and D. Pui, "Condensation-induced particle formation during vacuum pump down," *Journal of the Electrochemical Society*, vol. 140, pp. 1463–1468, 1993.

K. J. Budde, W. J. Holzappel, and M. M. Beyer, "Application of ion mobility spectrometry to semiconductor technology: Outgassings of advanced polymers under thermal stress," *Journal of The Electrochemical Society*, vol. 142, no. 3, pp. 888–897, 1995.

C.J. Clark, V.A.McGlone, A.White, A.B.Woolf, Dry matter determination in 'Hass' avocado by NIR spectroscopy *Postharvest Biology and Technology* , Volume 29, Issue 3, September 2003, Pages 301-308

P. R. Griffiths and J. A. De Haseth, *Fourier Transform Infrared Spectrometry*. Wiley, 2007.

T. Iwamoto and T. Ohmi, "Ultra-thin gate oxide reliability enhanced by carbon contamination free process," *Applied Surface Science*, vol. 117-118, pp. 237 – 240, 1997.

Jia-Huan Qu, Dan Liu, Jun-Hu Cheng, Da-Wen Sun, Ji Ma, Hongbin Pu &Xin-An Zeng (2015) Applications of Near-infrared Spectroscopy in Food Safety Evaluation and Control: A Review of Recent Research Advances, *Critical Reviews in Food Science and Nutrition*, 55:13, 1939-1954, DOI: 10.1080/10408398.2013.871693

Y. Kawai, A. Otaka, A. Tanaka, and T. Matsuda, "The effect of an organic base in chemically amplified resist on patterning characteristics using KrF lithography," *Japanese Journal of Applied Physics*, vol. 33, pp. 7023–7027, dec 1994.

J. A. Lebens, W. C. McColgin, J. B. Russell, E. J. Mori, and L. W. Shive, "Unintentional doping of wafers due to organophosphates in the clean room ambient," *Journal of The Electrochemical Society*, vol. 143, no. 9, pp. 2906–2909, 1996.

K. Saga and T. Hattori, "Influence of surface organic contamination on the incubation time in low-pressure chemical vapor deposition of silicon nitride on silicon substrates," *MRS Proceedings*, vol. 477, p. 379, 1997.

K. Siefering, H. Berger, and W. Whitlock, "Quantitative analysis of contaminants in ultrapure gases at the parts-per-trillion level using atmospheric pressure ionization mass spectroscopy," *Journal of Vacuum Science & Technology A*, vol. 11, no. 4, pp. 1593–1597, 1993.

Xiaonan Lu, 2017, *Sensing Techniques for Food Safety and Quality Control*, ISBN, Sensing Techniques for Food Safety and Quality Control, the Royal Society of Chemistry, ISBN 978-1-78262-664-0

Xingbo Xi et al, 2019, Review on carbon dots in food safety applications *Talanta* Volume 194, 1 March 2019, Pages 809-821

.

.

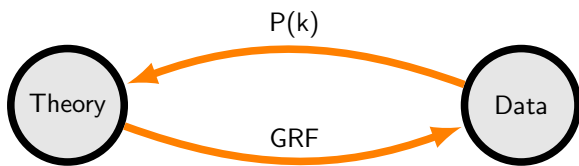
Numerical and Statistics methods for Cosmology

Fromenteau Sebastien

October 21, 2019

1 Introduction

The Gaussian Random Field is a powerful tool to produce random realization of random process which follow a Gaussian probability distribution function (PDF hereafter). Due to the Central Limit Theorem, random variable resulting from a sum of various different processes tend to follow a Normal distribution. We will see that generate a Gaussian Random Field is equivalent to generate random values following a Gaussian distribution (which will depend on the power spectrum) that will be the Fourier coefficients associated with the Fourier base function. We will then first introduce the Fourier Transform (i.e Fourier space decomposition). Moreover, we will see that during the early Universe (before the formation of the galaxies) the perturbations are small enough to consider independent evolution in time of the modes. For this reason, the Fourier space is of great importance for Cosmology. The Idea of this lecture is to understand how to generate a Gaussian Random Field following a given power spectrum. So, in a simple picture, the power spectrum estimation allows to link data to theoretical predictions, and the Gaussian Random Field (GRF) allows to generate data from a given power spectrum, mimicking the 2-pt statistics from a given theory.



The power spectrum $P(\vec{k})$ is the variance of the random process governing the value of $\delta_{\vec{k}}$ which is the Fourier coefficient associated to \vec{k} for a given real space field $\delta(\vec{r})$:

$$\delta(\vec{r}) = \sum_{\vec{k}} \delta_{\vec{k}} e^{i\vec{k} \cdot \vec{r}}. \quad (1)$$

In this equation, we do not specify the normalization which depends on the different existing conventions, however it do not impact the idea. We also have to keep in mind that the $\delta_{\vec{k}}$ are generally complex numbers. We will develop this expression for 1 dimension and then we will generalize to 2 and 3 dimensions.

2 Correlation Function and Power Spectrum

The Probability Distribution Function of a stochastic process contains the whole information about it. If you know this function you can derive all the quantities you want in order to compare them with data. However, most of the time we do not have access to this valuable information and we need estimate it. The correlation function is a powerful tool to do it. We can evaluate the N-pt correlation function in order to access to the different moment of the underlying probability distribution. On the other hand, we can also use their corresponding Fourier transform (Power Spectrum for the 2-pt CF, Bispectrum for the 3-pt CF, Trispectrum for the 4-pt CF...)

We will start to introduce the 2pt-correlation function which is more intuitive than the power spectrum at first look but we will more focus on the later for the gaussian Random Field use.

2.1 2pt Correlation Function

A way to characterize a random process is to estimate the correlation functions (different orders will correspond to the number of points we use). The simplest and more important is the 2-pt correlation function. If we have a random distribution of points with a mean density $\bar{n} = N/V$; where N is the number of points in the total volume V ; the number of points we can expect to see in a small volume dV is given by : $dN = \bar{n}dV$. If we consider two different small volume dV_1 and dV_2 we expect to observe in average $dN_1 = \bar{n}dV_1$ and $dN_2 = \bar{n}dV_2$ points respectively. So, we expect to get an average number of pairs of points between the two volumes equal to:

$$dP_{1,2} = dN_1 \times dN_2 = \bar{n}^2 dV_1 dV_2. \quad (2)$$

Because we generally use only the counting around existing points, we are interested in the number of pairs existing between this point and the points inside a small volume dV . It directly corresponds to the number of expected points in the later volume:

$$dP = 1 \times dN = \bar{n}dV. \quad (3)$$

If the distribution is not exactly random in positions, we expect to obtain a different quantity. So we let the

possibility to have an excess or default respect to the randomly expected value writing:

$$dP(\vec{r}) = 1 \times dN = \bar{n}dV[1 + \xi(\vec{r})], \quad (4)$$

where $\xi(\vec{r})$ is the 2-pt correlation function. If the distribution of the points is random, the number of pairs will be compatible with $\bar{n}dV$ and so the correlation function will be null. In the other case, we will find excess and defaults in particular directions and orientations. If we consider an isotropic distribution (like the Universe if we believe in the Cosmological Principle) the deviation from the random expected number of pair have to be independent of the orientation and so will depend only on the distance $|\vec{r}|$. Moreover, in the case isotropy, we can directly consider the all shell over the point with radius r . So we can recast the equation using the volume in the shell as:

$$dP(r) = \bar{n} \times 4\pi r^2 dr [1 + \xi(r)], \quad (5)$$

which is the most common way to express the 2pt-correlation function in cosmology. In all the reasoning we done before we use the number of pair of points we expect to measure so the natural way to estimate the correlation function will be using the pair counts at each scale r and compare it with the expected value for a random distribution. In the case of a simple realization, we can estimate the correlation function as:

$$\hat{\xi}(r) = \frac{DD(r)}{N\bar{n} \times 4\pi r^2 dr} - 1, \quad (6)$$

where $DD(r)$ is the number of pairs we count at a distance $\in [r, r + dr]$ considering the N Data points. As simple case, we refer to a periodic box for which there is no limits in the pair counting. Indeed, we can always draw a complete shell around each point in the limit of the size of this box. However, the reality is different and the expected number of pairs for the random realization is in general impossible to evaluate theoretically. For this reason, we create a random catalog reproducing the geometry containing the data points and we can then compare the pair counts between the data $DD(r)$ and the random realization $RR(r)$ as:

$$\hat{\xi}(r) = \frac{N \times DD(r)}{N \times RR(r)} - 1 = \frac{DD(r)}{RR(r)} - 1. \quad (7)$$

Moreover, we in order to reduce the variance in the estimation of $\hat{\xi}(r)$, we can increase the density of the random sample. We also need to take in to account that we will measure more pairs in the random than in the data points. If we multiply the density by a factor β then we will have $\beta N(\beta N - 1) \approx \beta^2 N^2$ pairs when we will measure $N(N - 1) \approx N^2$ pairs for the data:

$$\hat{\xi}(r) = \beta^2 \frac{DD(r)}{RR(r)} - 1, \quad \beta \approx \frac{\bar{n}_{random}}{\bar{n}_{data}}. \quad (8)$$

It exists different estimators to evaluate the 2pt-correlation function and the one optimizing the variance

and the bias of the estimation is provided by ?:

$$\hat{\xi}_{L-S}(r) = \frac{DD(r) - 2DR(r) + RR(r)}{RR(r)}, \quad (9)$$

where $DR(r)$ is the number of pairs we can do between points from the random and the points from the data. It is possible to do it since we reproduce the geometry of the data in the random catalog. While the form looks very simple, the demonstration to show the efficiency of this estimator is pretty hard and developed in the reference.

We will not develop in more details the 2-pt correlation function since we have enough material to understand the Gaussian Random Field.

2.2 Power Spectrum

In this section we will present how to calculate a power spectrum from 3 dimension statistically isotropic data. We used the power spectrum for generating the GRF without explain how to measure/estimate it. We will also briefly see how to obtain a theoretical prediction for the power spectrum in order to have an idea why we can use it as a link between theories and data.

The power spectrum contains the information over the variance of a random process. We then understand that the power spectrum exists independently of a Gaussian process, but it contains the all information only for the later case. Knowing that the inflation scenarios still allowed by the CMB constraints predicts a Gaussian or almost Gaussian field, we understand the importance of the power spectrum for the perturbations in cosmology. Moreover, we know from the CMB observations that these fluctuations are lower than 10^{-4} up to the recombination making negligible the cross mode terms contribution. It follows that we can evolve the k-mode perturbations (and so the associated wave-plane in real space) individually during the radiation era. It is what we call the linear perturbation theory.

2.2.1 Definition and Fourier Transform

The power spectrum contains the information of the variance for each corresponding Fourier k-mode. In case we have access to various realization of a same stochastic process we can calculate the Fourier Transform of each realization and obtain the distribution of the associated coefficients. Before to go further we need to briefly define the Fourier Transform and particularly the discrete transformation which is the only one we can compute on non trivial data.

Discrete Fourier Transform The Fourier space is the one formed by the sine and cosine functions, or in terms of exponential functions using the complex notations as presented in Eq.(1). We will use the later along this note. This basis is of first importance being naturally invariant under rotation and translation. We understand that under the Cosmological Principle, the

information will be optimally compress in the Fourier space.

2.2.2 Ergodicity and Cosmological Principle

As all statistical object, we need to estimate it from data. In the 1D and anisotropic cases (where we do not have redundant information in a stand alone realization) we need to work with several realizations or to use the ergodicity of the data. The ergodicity is the assumption of that various parts of the data are independent realizations of the same random process. In case of the 1D temporal case, it is equivalent to consider that the value $x(t)$ for each time t follows the same distribution allowing to infer the underlying probability distribution function (PDF) measuring the different moments over a long time serie of the data. In case of non-ergodicity we need to have access to different realizations $x_i(t)$ of the same random process and infer the different moments of the PDF at each time. In order to understand this point let us define properly the centered moments as:

$$\mu_1 = \int_{-\infty}^{\infty} x f(x) dx, \quad (10)$$

$$\mu_n = \int_{-\infty}^{\infty} (x - \mu_1)^n f(x) dx \quad \text{for } n > 1. \quad (11)$$

The way to determine them is to create an estimator over data expecting that the used sample will be representative enough of the underlying PDF $f(x)$. The estimators for the first two moments are given by:

$$\tilde{\mu}_1 = \frac{1}{N} \sum_{i=1}^N X_i; \quad \tilde{\mu}_2 = \frac{1}{N-1} \sum_{i=1}^N (X_i - \tilde{\mu}_1)^2 \quad (12)$$

where the pre-factor for $\tilde{\mu}_2$ is obtain using Jacknife calculation and show that the creation of an estimator is highly non-trivial. However, this topic is not the subject of this note and we will not discuss it. We will focus on the data we use to make the estimation. If we do not have any idea of the underlying process which generate the data, we should use different realization sets and do the summations over the number of these realizations. So N should be the number of the realizations and we should have as much estimations of $\tilde{\mu}_1$ and $\tilde{\mu}_2$ as number of points in each realizations. In the real case, we almost never have access to independent realizations and we use the data of one realizations as independent realizations of the underlying process. This assumption is called the ergodicity. So for the 1D time series, it is equivalent to do the summation over the time using just one realization (one time series). Considering now the 2D and 3D spatial cases, evaluate the mean density of a random field is done doing the summation over the different coordinates. However, when the CP stipulate that the Universe is statistically hogeneous it implies that the mean density measured over various realizations of universe is independent of the position \vec{r} :

$$\bar{\rho} = \bar{\rho}(\vec{r}) = \langle \rho(\vec{r}) \rangle_{\Omega}, \quad \Omega \text{ being the ensemble} \quad (13)$$

of universes realizations.

In cosmology, the ergodicity is tacitly used most of the time and one can be confused at the time to think about the Cosmological Principle. This principle postulate that "the Universe is statistically homogeneous and isotropic" while we can often listen that "the Universe is homogeneous and isotropic". The later proposition is obviously wrong if we do not add any precision about some specific scales. Our existence or the fact that there are structures in the Universe are contrary to this proposition, The term "statistically" is of great importance and stipulate that our Universe is only one realization of stochastic process and that if we have access to different universes resulting of the same process then the estimation of the mean density over the all realizations should be independent of the position.

2.2.3 Estimation of the power spectrum

We introduced all the concept and we can focus on the Power Spectrum definition, which is the Fourier transform of the 2pt correlation function. Another way to write the correlation between two random vectors $\mathbf{X} = \{X_0, \dots, X_n\}$ and $\mathbf{Y} = \{Y_0, \dots, Y_n\}$ is in term of the covariance:

$$C(\mathbf{X}, \mathbf{Y}) = \langle (X_i - \bar{X}) \cdot (Y_i - \bar{Y}) \rangle, \quad (14)$$

where the mean quantities vanish when considering contrast density variables. The power spectrum is the Fourier transform of this simple object and reads:

$$P(\vec{k}) = FFT[\xi(\vec{r})] = FFT[\langle \delta(\vec{x}) \cdot \delta(\vec{x} + \vec{r}) \rangle_{\vec{x}}],$$

$$P(\vec{k}) = \frac{1}{(2\pi)^3} \langle \hat{\delta}(\vec{k}) \cdot \hat{\delta}(\vec{k}') \rangle_{\Omega} \delta^D(\vec{k} + \vec{k}'), \quad (15)$$

where $\hat{\delta}(\vec{k})$ is the Fourier transform of the contrast density field and $\delta^D(\vec{k} + \vec{k}')$ is the Delta Dirac function which appears to guaranty the homogeneity. Here, we need to have access to independent realizations of the $\hat{\delta}(\vec{k})$ which can be done dividing the volume of the data in "independent" sub volumes. So we can write simplify using the fact that the Delta Dirac function is non null only if its argument is the null vector (so when $\vec{k}' = -\vec{k}$):

$$P(\vec{k}) = \frac{1}{(2\pi)^3} \langle \hat{\delta}(\vec{k}) \cdot \hat{\delta}(-\vec{k}) \rangle_{\Omega}. \quad (16)$$

We will see in Sec 3.2 that in case of real field (i.e. not imaginary part), that $\hat{\delta}(-\vec{k}) = \hat{\delta}(\vec{k})^*$. It comes that the Power Spectrum for the vector \vec{k} reads:

$$P(\vec{k}) = \frac{1}{(2\pi)^3} \langle \hat{\delta}(\vec{k}) \cdot \hat{\delta}(\vec{k})^* \rangle_{\Omega} = \frac{1}{(2\pi)^3} \langle |\hat{\delta}(\vec{k})|^2 \rangle_{\Omega} \quad (17)$$

Moreover, the statistical isotropy also apply here. As we will see further, the $\hat{\delta}(\vec{k})$ are the amplitude and phases of the pane wave with frequency $|\vec{k}|$ and propagation vector $\vec{k}/|\vec{k}|$. The isotropy stipulate that the information

do not depend on the orientation which implies that the information bring by each vector \vec{k} such as $|\vec{k}| = k_1$ is an independent realization of the same stochastic process at the scale k_1 . Then, we can reduce the isotropic Power Spectrum definition as:

$$P(k) = \frac{1}{(2\pi)^3} \left\langle |\hat{\delta}(\vec{k})|^2 \right\rangle_{|\vec{k}|=k}$$

3 Gaussian Random Field

Now we introduce properly the Fourier Transform and the power spectrum we have all the tools in order to understand and generate Gaussian Random Fields. Contrary to the estimation of the power spectrum, there is no difference on the GRF generation considering isotropic and anisotropic cases. So no distinction is needed. We first present the 1D case before to extend to the 2D and 3D cases. As we will see, there is no main differences between the different dimensions except the visualization. However, the faculty to visualize is important to infer and interpret the results for the structure formation. The 1D case will help to accustom with the random generation. We can particularly use it to support the Press and Schechter formalism in the Halo Mass Function generation. Then, we will focus on the 2D case which corresponds to the Cosmic Microwave Background space. We will first work on the wave-plane approximation before to introduce the spherical decomposition using the Legendre polynomials. It will be helpful for the $C_\ell(X_1, X_2)$ cross correlation where X_1 and X_2 are the perturbations measurements from the CMB's photons for the temperature, the E-modes or B-modes of polarizations. Finally we will treat the 3D case which is particularly in relation with the Cosmological and Astrophysical (up to galaxy scales) simulations. Indeed, the initial conditions of a large scale simulations are relied to the primordial power spectrum or with the CMB's power spectrum.

3.1 1D-Gaussian Random Field

We start with the 1D case which allows to write simply the relations we need. First, we remind that we generally use a central random field (*i.e* the mean of this field is 0) reason why we will use the notation δ in all this lecture. For example, we use the density perturbation defined as $\delta_\rho = \frac{\rho - \bar{\rho}}{\bar{\rho}}$. The time series analysis in finance use a lot the "return" variable which corresponds to the relative variation of the price at each time step. In both cases, the mean of the variable is 0.

Anyway, in one dimension, we can rewrite Eq. 1 as

$$\delta(r) = \sum_{k=-k_{max}}^{k_{max}} \delta_k e^{ik.r}, \quad (18)$$

$$\delta(r) = \sum_{k=k_{min}}^{k_{max}} (\delta_k e^{ik.r} + \delta_{-k} e^{-ik.r}) + \delta_0 \underbrace{e^{i0.r}}_{=1} \quad (19)$$

where we regroup the terms in k and $-k$ together and where we exit the $k = 0$ term from the sum. We can see that the term $\delta_{k=0}$ contribute as a constant and so will just shift the field $\delta(r)$ or can be interpreted as the mean of the density field. If we use a contrast density field variable, *i.e* $\delta = \frac{x - \bar{x}}{\bar{x}}$ where x is a stochastic variable, the mean is null by definition.

3.2 Condition for real field (no imaginary part)

In general, we want to generate a real gaussian random field in physics.....because most of the natural field are reals. So we can look for the general condition which guaranty that $\delta(r)$ is real. We will start from the result and verify that it always allows the reality of $\delta(r)$. This condition is :

$$\delta_{-k} = \delta_k^*. \quad (20)$$

Indeed, restarting from Eq.18 and including the condition Eq. 20 we get

$$\delta(r) = \sum_{k=k_{min}}^{k_{max}} (\delta_k e^{ik.r} + \delta_k^* e^{-ik.r}) + \delta_0, \quad (21)$$

$$\delta(r) = \sum_{k=k_{min}}^{k_{max}} \left(\underbrace{\delta_k e^{ik.r} + [\delta_k e^{ik.r}]^*}_{=2 \times \text{Re}(\delta_k e^{ik.r})} \right) + \delta_0, \quad (22)$$

where we can see that all the terms are purely reals if we fix the shift $\delta_{k=0}$ to be real.

3.2.1 Random generation of the δ_k

Now we define the condition for real field, we can generate a Gaussian Random Field which corresponds to generate the δ_k values. In order to be more explicit, we can write the complex number δ_k with its module α_k and phase ϕ_k :

$$\delta_k = \alpha_k \times e^{i\phi_k} \Rightarrow \delta_k^* = \alpha_k \times e^{-i\phi_k}. \quad (23)$$

We can see that if we only generate the modules α_k letting the phases $\phi_k = 0$ we get a specific realization of a real gaussian random field.

Generate α_k The most important part of the Gaussian Random Field is contained inside the α_k . Indeed, the word "Gaussian" is associated with the probability distribution function of these variables. And because the δ_k are variable of zero mean, the only information we need is the standard deviation or the variance. As we seen during the last lecture, the variance of the process is provided by the power spectrum $P(k)$, so generate a α_k corresponds to generate a random value following the pdf:

$$\alpha_k \sim \mathcal{N}(0, \sigma^2 = P(k)).$$

The word "Random" comes from this random generation of the α_k following a "Gaussian" probability distribution with variance provided by the power spectrum.

Generate ϕ_k In general, we also need to generate the phases. However, in cosmology we do not have (at least we think that there is not) information from the phase. So we also generate the phases using a uniform distribution between 0 and 2π . One time we get a phase we have to remember to satisfy the condition Eq.20 to produce a real random field.

$$\phi_k \sim \mathcal{U}([0, 2\pi]) \quad \& \quad \phi_{-k} = -\phi_k.$$

3.2.2 Get the random field $\delta(r)$

Finally, we have to calculate the inverse Fourier transformation Eq.1 to get the real space gaussian random field. We will show basic results for various cases in order to well understand in details the different impacts we describe above.

3.2.3 Individual mode impact

In order to understand the real space gaussian random field let see the impact of the individual modes in a 1D

$$\begin{aligned} \delta(\vec{r}) &= \sum_{k_x=-k_{x,max}}^{k_{x,max}} \sum_{k_y=-k_{y,max}}^{k_{y,max}} \sum_{k_z=k_{z,min}}^{k_{z,max}} \left(\delta_{\vec{k}=(k_x, k_y, k_z)} e^{i(k_x \cdot r_x + k_y \cdot r_y + k_z \cdot r_z)} + \delta_{-\vec{k}=(-k_x, -k_y, -k_z)} e^{-i(k_x \cdot r_x + k_y \cdot r_y + k_z \cdot r_z)} \right), \\ \delta(\vec{r}) &= \sum_{k_x=-k_{x,max}}^{k_{x,max}} \sum_{k_y=-k_{y,max}}^{k_{y,max}} \sum_{k_z=k_{z,min}}^{k_{z,max}} \left(\delta_{\vec{k}=(k_x, k_y, k_z)} e^{i\vec{k} \cdot \vec{r}} + \delta_{-\vec{k}=(-k_x, -k_y, -k_z)} e^{-i\vec{k} \cdot \vec{r}} \right), \\ \delta(\vec{r}) &= \sum_{k_x=-k_{x,max}}^{k_{x,max}} \sum_{k_y=-k_{y,max}}^{k_{y,max}} \sum_{k_z=k_{z,min}}^{k_{z,max}} \left(\delta_{\vec{k}=(k_x, k_y, k_z)} e^{i\vec{k} \cdot \vec{r}} + \left[\delta_{\vec{k}=(k_x, k_y, k_z)} e^{i\vec{k} \cdot \vec{r}} \right]^* \right), \\ \delta(\vec{r}) &= \sum_{k_x=-k_{x,max}}^{k_{x,max}} \sum_{k_y=-k_{y,max}}^{k_{y,max}} \sum_{k_z=k_{z,min}}^{k_{z,max}} \left(\underbrace{\delta_{\vec{k}} e^{i\vec{k} \cdot \vec{r}} + \left[\delta_{\vec{k}} e^{i\vec{k} \cdot \vec{r}} \right]^*}_{=2 \times Re(\delta_{\vec{k}} e^{i\vec{k} \cdot \vec{r}})} \right). \end{aligned} \tag{24}$$

So we get the same result than for the 1D case and the general condition for real Gaussian random field reads:

$$\text{Reality} \iff \delta_{-\vec{k}} = \delta_{\vec{k}}^*$$

3.3.2 Isotropic field

Moreover the Cosmological Principle postulate that the Universe is statistically Isotropic and Homogeneous.

case. Let start with a simple example in which we use a power spectrum in $P(k) = k^{-2}$. We generate the values of the δ_k , respecting the reality condition, that we show the result for various k modes in the figure 3. Because we give more importance to the lower k values, we conserve a general form at large scale. The details at small scales (*i.e* high k values) are imprinted continuously but with diminishing amplitudes. The figure shows the individual contribution of each mode associated with the sum of all the contributions from the lower modes (*i.e* $\sum_{k' \leq k} \delta_{k'} e^{ik'x}$).

3.3 2D and 3D Gaussian Random Field

One time we understand how works the generation in 1D, the generalization to 2D and 3D is relatively straightforward. We will have to specify the isotropic/anisotropic conditions and their implications on the δ_k generation. In case of isotropy, we will see the reason why we can define a power spectrum $P(k)$ which depends only on the module of \vec{k} and so understand why we generally we use it in cosmology. We will mostly illustrate our purpose in 2D for graphic convenience but the results are totally similar in 3D.

3.3.1 Reality of the Gaussian Random Field in 2D and 3D

We have seen in section 3.2 that the condition for reality is $\delta_{-k} = \delta_k^*$ for 1D scalar k . We will demonstrate now that the general condition of reality is $\delta_{-\vec{k}} = \delta_{\vec{k}}^*$ independently of the dimension. So we will do the demonstration only for 3D which embed the 2D case.

We can think quickly about the signification using the basic definition Eq.1 and see that if the information depends on the orientation of \vec{k} corresponds to have

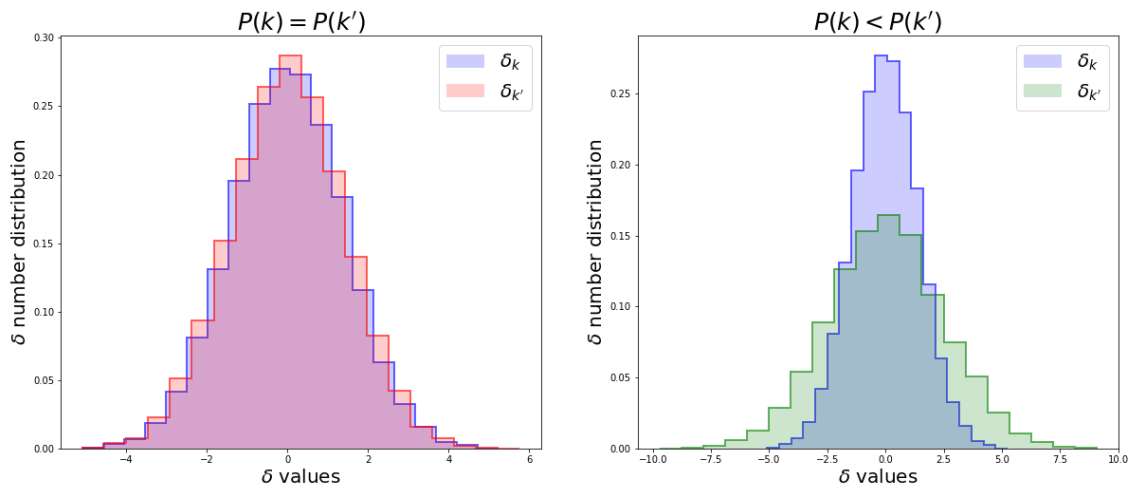


Figure 1: *Left panel* : distribution of 10000 realizations for two δ_k using the same power spectrum value. *Right panel* : distribution of 10000 realizations for two δ_k using different power spectrum values. We can see that a larger power spectrum value allows larger absolute values for the δ .

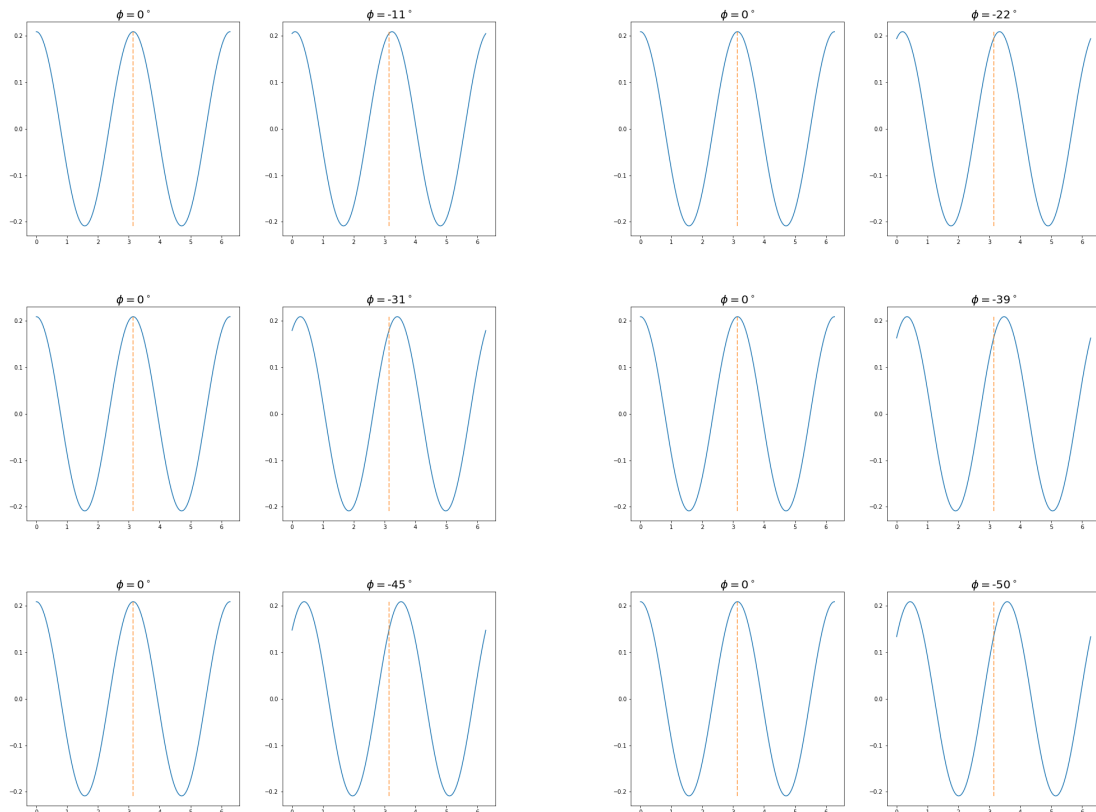


Figure 2: Impact of the phase ϕ on a given mode in real space. We can see that it lies a shift on the cosine wave.



Figure 3: Impact of the individual modes on a 1D Gaussian random field following a power spectrum in k^{-2} . Each one of the 10 plots shows the contribution of a specific mode (the right part) and the sum of all the contributions of the modes lower and equal to this mode.

a probability distributions of the $\delta_{\vec{k}}$ depending on the orientation of \vec{k} . As we seen in section 3.2.1, it corresponds to give in average more amplitude to the modes oriented in direction of the \vec{k} vectors with greater values in the power spectrum $P(\vec{k})$. We will express it with more details. We consider two vectors $\vec{k} = (k_x, k_y)$ and $\vec{k}' = (k'_x, k'_y)$ with same module $|\vec{k}| = |\vec{k}'|$ and we will compare the impact on the real space.

If we have the same variance for the two modes (*i.e* $P(\vec{k}) = P(\vec{k}')$), then we have that the two probability distributions are equals:

$$\delta_{\vec{k}}, \delta_{\vec{k}'} \sim \mathcal{N}\left(0, \sigma^2 = P(\vec{k}) = P(\vec{k}')\right), \quad (25)$$

and so we expect to have over a large enough sample a similar distribution for two variables (Left panel in figure 1). In the figure 4, we can see the contribution in the real space field of various vectors \vec{k} with same module. As we can expect considering the dot product $\vec{k} \cdot \vec{r}$, the plane wave is oriented along \vec{k} . So, if we give more importance to the mode \vec{k}' than to the mode \vec{k} (*i.e* $P(\vec{k}) < P(\vec{k}')$), then the amplitude of the wave oriented along \vec{k}' will be in general more important than the one oriented along \vec{k} and generate an anisotropy in the real space. So, it appears that the condition to produce an isotropic Gaussian Random Field we need the condition:

$$\text{Isotropy} \iff P(\vec{k}) = P(|\vec{k}|)$$

3.3.3 Impact of individual modes in 2D and 3D

The 2D Fourier decomposition corresponds to a sum of plane wave with constant values perpendicular to the vector \vec{k} as we can see an example in figure 5. The reason is simple, along an axe perpendicular to \vec{k} , the dot product $\vec{k} \cdot \vec{r}$ is constant, so the term in $e^{i\vec{k} \cdot \vec{r}}$ is constant to. In 3D, the constant value will be for all \vec{r} reaching the plane perpendicular to \vec{k} as shown in figure 6.

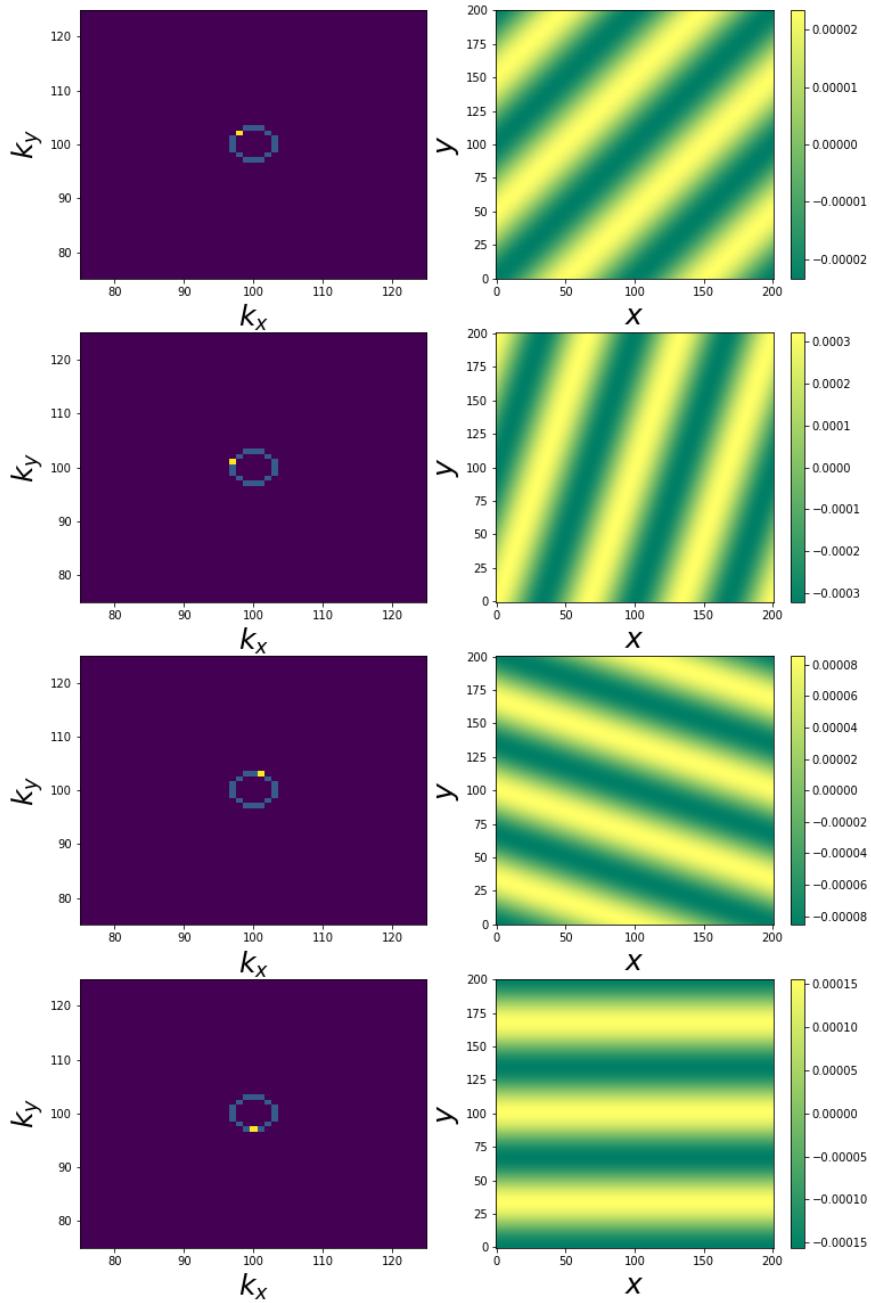


Figure 4: Effect in real space of individual pixels with same module. Same module $|\vec{k}|$ implies that the frequency/wavelength is the same for all the wave plane in the plots. The left panels show the pixel in the Fourier space which was used to generate the wave plane in the right panels. We can see that the iso-values in the right panels are perpendicular to the Fourier vector \vec{k} as explained in the text. The amplitude of the wave plane fluctuations are different because the $\delta_{\vec{k}}$ differs following a random trial.

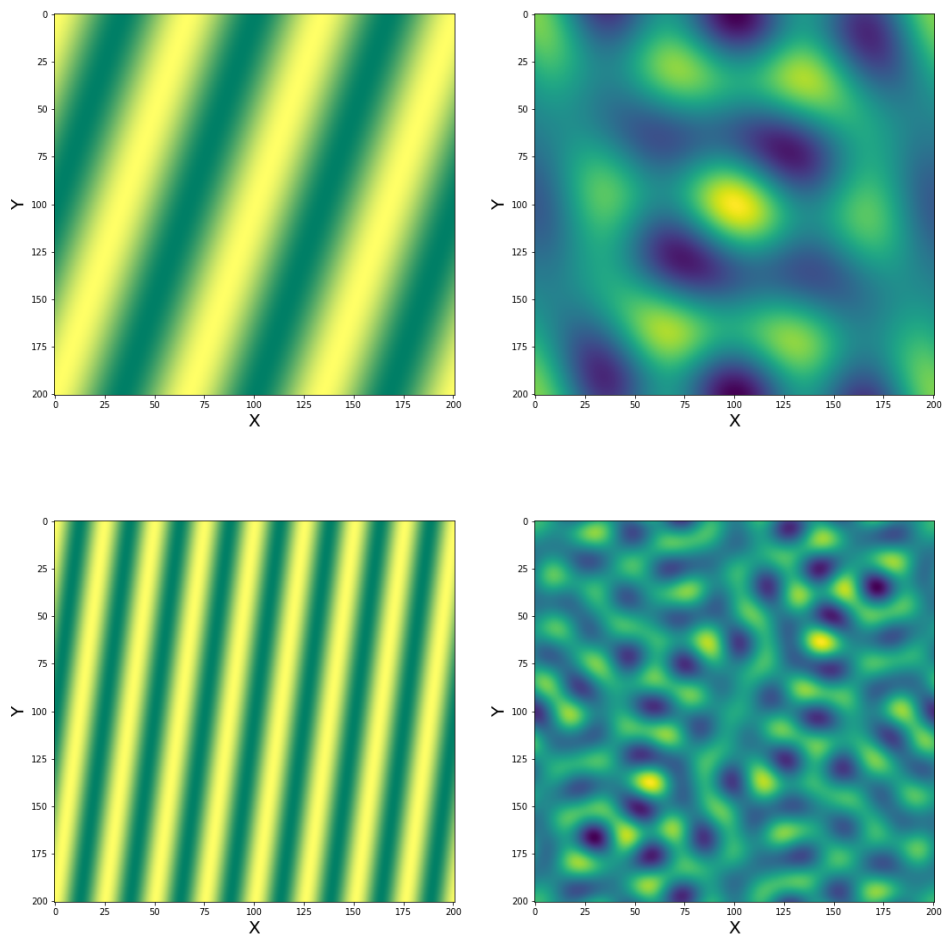


Figure 5: *Top Left panel:* Real space result for 1 pixel in Fourier space. *Top Right:* Real space result of the sum of the pixel in Fourier space with same module. The pixels used are all points which compose the circle we can see on the left panels of the figure 4. *Bottom Left and Right:* We redo the same exercise for another frequency, so another module of \vec{k} which is greater than the previous one.

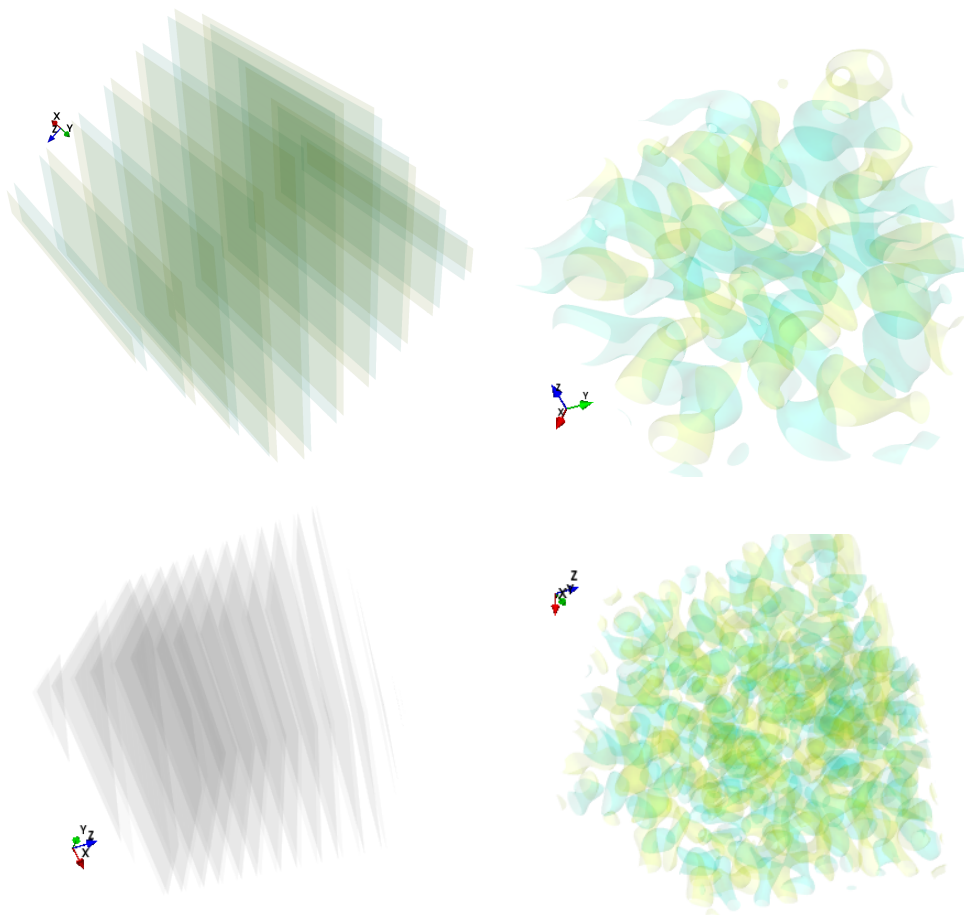


Figure 6: *Top Left panel* : Real space result for 1 pixel \vec{k} in 3D. The result shows the iso-contours which corresponds to perpendicular planes as explained in the text. Each plane corresponds to the same value so the space between 2 planes correspond to the wavelength of the mode \vec{k} . *Top Right panel* : Real space result of the sum of all the pixels on a sphere (a shell) so the contribution of all the \vec{k}_i tq $|\vec{k}_i| = k \forall i$. *Bottom Left and Right panels* : The same than above but for a greater module $|\vec{k}|$

LA FÍSICA FUNDAMENTAL Y APLICACIONES DE LOS IONES NEGATIVOS

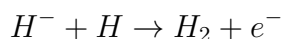
Guillermo Hinojosa Aguirre

Instituto de Ciencias Físicas
de la Universidad Nacional Autónoma de México,
UNAM, Campus Morelos.

INTRODUCCIÓN

Desde que Wild descubrió, en los años treinta, que el H^- es el portador de la opacidad interestelar en 1650 nm y no los metales como se pensaba [1], los iones negativos han sido un pregunta fundamental tanto por su rareza estructural como por su gran abundancia en plasmas y en ambientes atmosféricos.

Desde entonces, los iones negativos se usan para explicar diversos fenómenos en más de una campo de la ciencia. Por ejemplo, en astrofísica se considera que en la era de la recombinación (4×10^5 años después del *Big Bang*) el H^- contribuyó con el mecanismo para la formación de H_2 via la reacción



Una vez formado, el H_2 puede reducir la presión de forma eficiente en un orden de magnitud debido a sus transiciones roto vibracionales. Esto hizo que el tiempo de enfriamiento fuera menor que el de expansión y entonces se presentaron las condiciones para formar objetos estelares. Este es el mecanismo por el que los astrofísicos creen que gas interestelar a 1000 K pudo haber conducido a la formación de las primeras estrellas (de 10^6 masas solares) [2].

Más recientemente, la presencia de iones negativos en diversos ambientes atmosféricos y también en el espacio interestelar se ha confirmado. Estos iones tienen abundancias suficientes para detectarlos. Por ejemplo, en el coma del cometa Halley [3], en la región interestelar de Cepheus y Auriga [4] y particularmente, en la atmósfera de Titán [5].

La presencia de iones negativos en el espacio es contraintuitiva. En el espacio, se espera que la radiación destruya a los aniones debido a que el

electrón extra se encuentra muy débilmente ligado [6]. Por lo tanto, debe existir un mecanismo eficiente de formación de estos iones capaz de mantener las poblaciones en las cantidades que se han detectado. El en caso de los plasmas, la misma razón debería explicar bajas poblaciones de aniones debido a que con las interacciones el electrón se pierde fácilmente. Sin embargo, poblaciones altas de iones negativos son muy comunes en plasmas. Además, los aniones en plasmas están correlacionados con la presencia de polvo. El mecanismo por el que se forma el polvo en estos ambientes es desconocido, sin embargo se sabe que, en plasmas, los iones negativos y el polvo van de la mano [7].

El problema con la detección de esta clase de iones en estos ambientes hostiles es que debería de haber un mecanismo relativamente simple y eficiente para su formación, tal como la captura radiativa (CR) que es el mecanismo más canónico para la explicación de su formación,



en donde un electrón e^{-} se acerca lo suficiente para el el átomo o molécula M lo capture seguido de la emisión de radiación $h\nu$.

Sin embargo, la evidencia de laboratorio parece indicar que este mecanismo no es eficiente para justificar la cantidad de aniones observados [8]. Como consecuencia, la pregunta de cómo se forman los iones negativos en ambientes como el de los plasmas o del espacio interestelar sigue sin una respuesta completa y representa un reto para la ciencia actual.

INTERACCIONES DE ANIONES MOLECULARES CON GASES ATMOSFERICOS

Con la curiosidad de tratar de entender más sobre esta clase rara de iones, en el laboratorio hemos abordado el estudio de sus interacciones con gases comunes en ambientes atmosféricos como el nitrógeno (N_2) y el oxígeno (O_2). La idea es que a través de la medición de secciones transversales de despojo electrónico (el proceso inverso al de su formación) y más recientemente, al estudio de su fragmentación (con alta resolución) es posible conocer indirectamente algunos de las procesos de formación o de algunas de sus características (al menos esta es la hipótesis).

Además, las secciones transversales de despojo electrónico (σ) son de interés en el estudio del equilibrio de plasmas. A energías altas, la física de estas interacciones se puede explicar con relativo éxito mediante el modelo de esfera dura para colisiones. Bajo esta perspectiva, todas las secciones totales de

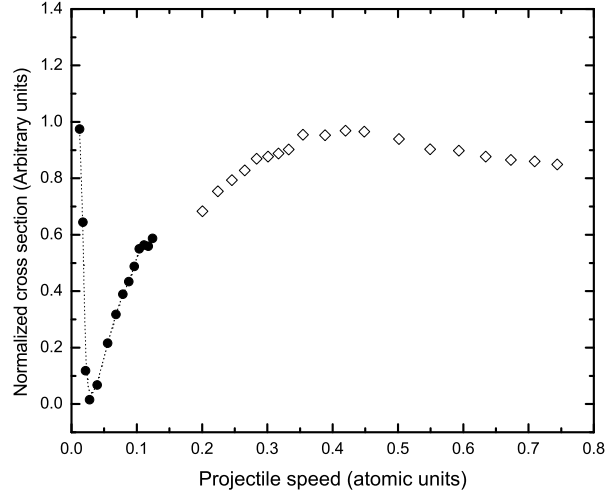


Figura 1: Gráfica tomada de [10]. Secciones transversales totales de despojo electrónico σ normalizadas según el modelo de colisión de esfera dura (con el factor de la ecuación 2). Los valores de σ representados con un diamante corresponden a Br^- en O_2 y con un círculo a $\text{CCH}_2^- + \text{O}_2$.

despojo electrónico siguen el mismo comportamiento cuando se normalizan por el área efectiva de interacción, que es el área total que subtienden las áreas del proyectil y del gas objetivo o blanco [9]. Esta es la clase de generalizaciones que nos gustan a los físicos. Es decir, si se conoce una sola sección (digamos σ_k) entonces los valores derivados de

$$\frac{\sigma_k}{\pi(r_p + r_t)^2}, \quad (2)$$

caen en la misma curva en donde todas las demás secciones para el mismo proceso también caen; r_p es el radio de los iones del proyectil y r_t es el radio de los átomos o moléculas del gas objetivo o blanco. Esto se ilustra en la Fig. 1.

Como consecuencia de esta hipótesis, si se conocen los radios de los nuevos iones proyectil y objetivo r_{po} y r_{to} entonces, cualquier otra sección σ_n se puede derivar simplemente de

$$\sigma_n = [r_{po} + r_{to}]^2 \frac{\sigma_k}{(r_p + r_t)^2} . \quad (3)$$

Sin embargo, encontramos que a energías bajas, σ no sigue este comportamiento [10, 11]. Esta misma tendencia, a bajas energías de interacción, lo acaba de confirmar un grupo europeo [12]. La razón por la que es importante este fenómeno es de origen práctico: los electrones de muy baja energía son de interés en varios campos de la física como en plasmas o en física médica.

Además encontramos un efecto adicional que consiste en la observación de estados metaestables de auto despojo que pensamos se forman durante la colisión y que tiene consecuencias directas en ciencia de plasmas. Resulta que este mecanismo de auto despojo ofrece una contribución de al menos un orden de magnitud mayor que el mecanismo de despojo por interacción de esfera dura.

Esto se logró a través de la medición de un parámetro de la física que indica la probabilidad por partícula en un área dada de que ocurra una cierta reacción (en este caso el despojo electrónico) después de interactuar con una molécula (en este caso N_2 u O_2). Este parámetro es mejor conocido como *Sección Transversal Total de Despojo Electrónico* y su símbolo es normalmente σ .

σ se midió con dos métodos experimentales. En uno de los métodos se miden los átomos neutros y en el otro se mide la atenuación de la intensidad del haz original de aniones. Ambos métodos miden el mismo proceso: la pérdida del electrón extra de los iones negativos. Ambos métodos tienen el mismo resultado final: un átomo o molécula neutra (que se monitorea experimentalmente) y un electrón libre (que no detectamos).

Para explicar esta técnica combinada hago uso de la Fig. 2 en donde se ilustra una región en donde el haz de aniones (que viene de la izquierda) interactúa con el gas de N_2 u O_2 . Como resultado de esta interacción la probabilidad de que un anión del haz pierda su electrón extra aumenta. Luego hay un campo eléctrico que desvía a los aniones remanentes del haz hacia un detector. Las partículas o átomos que pierden el electrón quedan neutras y el campo eléctrico no las puede desviar, siguen una trayectoria recta. Como consecuencia de incrementar la densidad del gas (o la presión) la cantidad de partículas neutras aumenta y la cantidad de haz de aniones disminuye.

El primer método se conoce como SGR (de sus siglas en inglés *signal growth rate method*) o método del crecimiento de señal y el segundo se conoce como BAT (de sus siglas en inglés *beam attenuation technique* o método de la

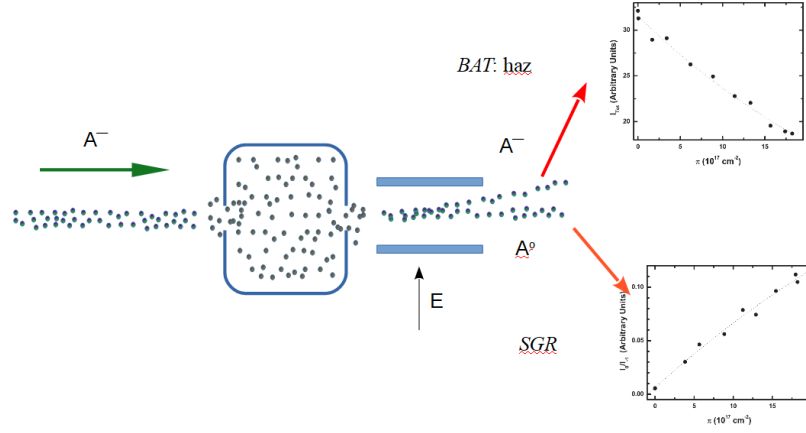
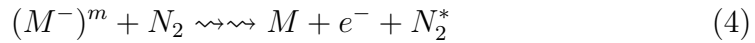


Figura 2: Se ilustra la idea de cada uno de los métodos SGR y BAT. A la zona en donde están confinados los gases se le conoce como celda de gas. El campo eléctrico entre las dos placas desvía la trayectoria de los aniones del haz hacia un detector situado fuera del eje central. Las partículas neutras que se forman entre la celda y las placas siguen sin desviarse hacia el detector central.

atenuación del haz. Ambos métodos miden el mismo fenómeno, la diferencia es que BAT mide la atenuación *total* del haz.

Lo que encontramos es que las σ medidas con BAT son mayores que las σ medidas con SGR (ver Fig. 3). Encontrar una explicación para estas diferencias fue un reto, pero después de mucha experimentación llegamos a la conclusión de que la diferencia tiene que ver con un mecanismo que genera más partículas neutras que no se toman en cuenta o no se detectan en SGR.

Aunque parezca increíble, se trata de una población de aniones del haz que interacciona con el gas y no pierde el electrón extra sino hasta que ha transcurrido un tiempo suficiente para que la partícula neutra resultante no sea detectada por el detector central. Eso quiere decir básicamente que el anión de haz pierde el electrón extra cuando va en camino al detector lateral. Dicho de otra forma, hay estados $(M^-)^m$ que se tardan (microsegundos) en perder el electrón extra y decaen lentamente,



En la Fig. 3 se muestra como para el caso de la molécula CCH_2^- (que es muy importante en plasmas de hidrocarburos) las σ son muy diferentes para

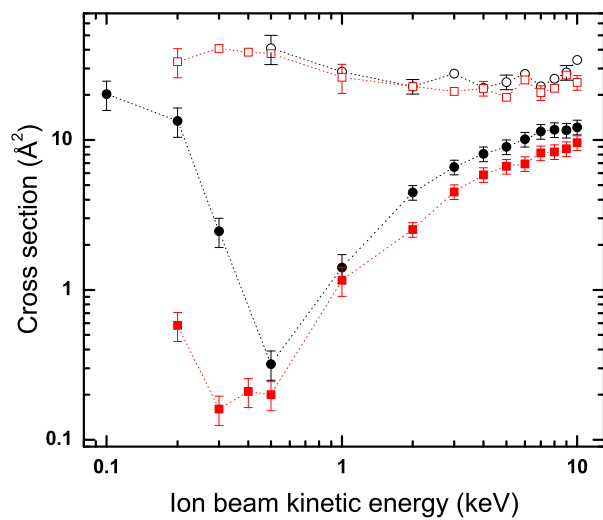


Figura 3: Gráfica tomada de [10]. Sección transversal de despojo electrónico inducida por colisión con gases atmosféricos de la molécula CCH_2^- . Los círculos corresponden a O_2 y los cuadros a N_2 . Los círculos y cuadros sólidos corresponden a σ^s y los huecos a σ^b .

energías de interacción bajas y como a energías altas las secciones convergen, sus valores se parecen más. Este es uno de los efectos que nos hicieron pensar en estados $(M^-)^m$.

La razón por la que estos estados de auto despojo no tienen efecto a energías altas es el tiempo de vida de estos estados metaestables. A energías altas el tiempo de vuelo de los aniones es muy corto y aparentemente no contribuyen a la sección transversal de despojo electrónico. La diferencia entre las secciones decae como función de la velocidad

$$\sigma^b - \sigma^s \quad (5)$$

en donde σ^b es la sección transversal de despojo electrónico medida con el método BAT y σ^s corresponde a la medición con SGR. Las partículas neutras que se generan como resultado de interactuar una sola vez con el gas de la celda vuelan dentro del experimento (que está a una presión muy baja). Este *tiempo de vuelo* depende solamente de la velocidad de haz ya que las dimensiones del experimento son constantes. Entonces a mayor velocidad de haz de aniones (o a energías más altas) el tiempo de vuelo de las partículas del haz es menor.

Si uno relaciona el tiempo de vuelo con el recíproco de la diferencia de la Ec. 5,

$$\beta = \frac{1}{\sigma^b - \sigma^s} \quad (6)$$

queda una gráfica como la de la Fig. 4. De estas gráficas propusimos derivar una especie de tiempo de vida (al que llamamos τ^β) que representa un valor del tiempo de vida de los estados metaestables $(M^-)^m$.

CONCLUSIONES Y PERSPECTIVAS

El proceso de despojo de electrones que sufren las moléculas simples como resultado de interacciones con gases atmosféricos no es tan simple como interacciones entre esferas duras. Es mucho más complejo. Un análisis basado en diferencias entre σ^b y σ^s reveló la presencia de estados metaestables de auto despojo que aparentemente se inducen durante la colisión. Estos estados $(M^-)^m$ tienen tiempos de vida del orden de μs y cuyo tiempo de vida τ^β medimos por primera vez. Es $4.4 \pm 4 \mu s$ para CCH_2^- en O_2 y $4.2 \pm 2 \mu s$ en N_2 . Es básicamente el mismo tiempo de vida. El tiempo de vida natural para el estado base de CCH_2^- es de 102 s.

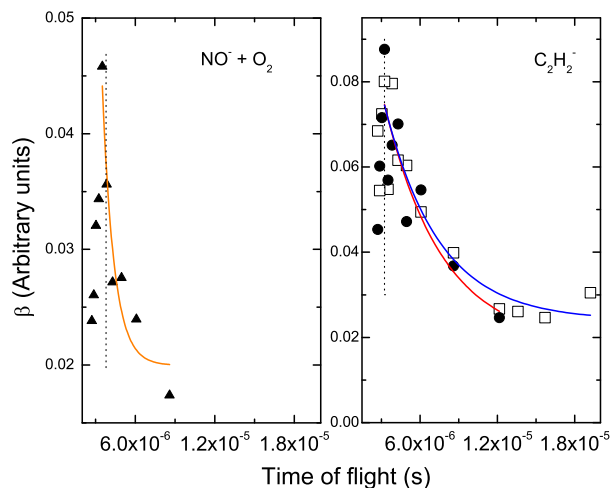


Figura 4: Decaimiento de la población de estados metaestables $(M^-)^m$ como función del tiempo de vuelo para dos especies moleculares. De aquí se derivaron tiempos de vida τ^β .

Finalmente, el anuncio que hice sobre un descubrimiento ya se confirmó y los resultados ya están publicados en la revista de física más reconocida en física. La misma revista en donde se acaban de publicar los resultados del LIGO. Es un hallazgo de carácter fundamental que va a cambiar la interpretación de los iones negativos y también la interpretación de todos los plasmas basados en hidrocarburos. Se trata de una especie nueva. De esto hablaré en la próxima escuela de verano en física ◁

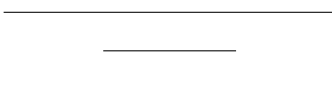
Referencias

- [1] R. Wildt. Negative ions of hydrogen and the opacity of stellar atmospheres. *Astrophys. J.*, 90:611–620, 1939.
- [2] T. J. Millar and et al. Negative ions in space. *Chem. Rev.*, 117:1765–1795, 2017.
- [3] Chaizy P. and et al. Negative ions in the coma of comet halley. *Nature*, 349:393, 1991.
- [4] Cordiner and et al. Discovery of interstellar anions in cepheus and auriga. *The Astrophysical Journal Letters*, 730:L18, 2011.
- [5] P. Lavvas and et al. Aerosol growth in titan’s ionosphere. *PNAS*, 110:2729–2734, 2013.

- [6] Cordiner M. A. and et al. On the ubiquity of molecular anions in the dense interstellar medium. *The Astrophysical Journal*, 770:48, 2013.
- [7] J. Winter. Dust in fusion devices a multi-faceted problem connecting high- and low-temperature plasma physics. *Plasma Phys. Control. Fusion*, 46:B583ÚB592, 2004.
- [8] M. Khamesian and et al. Formation of cn^- , $c3n^-$, and $c5n^-$ molecules by radiative electron attachment and their destruction by photodetachment. *Phys. Rev. Lett.*, 117:123001, 2016.
- [9] R. F. Nascimento and et al. Total detachment cross sections of c^- , c^- , $c2^-$, and $c2h^-$ incident on $n2$ at kev energies. *Phys. Rev A.*, 87:062704, 2013.
- [10] E. M. Hernández and G. Hinojosa. Collision induced electron detachment cross sections of the $h2cc^-$ anion below 10 kev on $o2$ and $n2$. *International Journal of Mass Spectrometry*, 424:35Ú39, 2018.
- [11] Serkovic-Loli and et al. Electron detachment of no^- in collisions with $o2$ and $n2$ below 10 kev. *Int. J. of Mass Spectrometry*, 392:23Ú27, 2015.
- [12] M. Mendes, C. Guerra, A. I. Lozano, D. Rojo, J. C. Oller, P. Limão Vieira, and G. García. Experimental electron-detachment cross sections for collisions of O_2^- with n_2 molecules in the energy range 50–7000 ev. *Phys. Rev. A*, 99:062709, 2019.

AGRADECIMIENTOS

El conocimiento nuevo que aquí se mostró fue financiado por la UNAM. Parte de los fondos corresponden al proyecto PAPIIT IN-109-317. Electrodo, analizadores, electrónica y software del experimento fueron desarrollados parcialmente por Guillermo G. Bustos M., Armando Bustos G., Héctor H. Hinojosa G., Reyes García C. y Arturo S. Quintero C.



Entrelazamiento, no-localidad y otras particularidades de la mecánica cuántica

F. Leyvraz

*Instituto de Ciencias Físicas—Universidad Nacional Autónoma de México, Cuernavaca, Morelos, México
and Centro Internacional de Ciencias, Cuernavaca, Morelos México**

(Dated: 15 de noviembre de 2019)

Se da una breve introducción a las ideas acerca de no-localidad en la mecánica cuántica y su relación con el entrelazamiento. Se da especial énfasis al concepto de *contrafactualidad*. Esto se refiere a nuestras creencias acerca de lo que hubiera pasado si algún determinado evento *hubiera ocurrido* de manera distinta a la en que efectivamente sucedió. Argumentaré que dichas creencias sirven para caracterizar una teoría clásica. Por otro lado, mostraré que, si consideramos tales hipótesis contrafactuales, la localidad efectivamente no se cumple en la mecánica cuántica. Se mostrará usando dos ejemplos bien conocidos, la “paradoja” debida a Lucien Hardy así como las correlaciones que se observan en el llamado estado de Greenberger, Horne y Zeilinger, que es un estado involucrando el entrelazamiento entre 3 espines.

I. CARACTERÍSTICAS DE LA MECÁNICA CUÁNTICA

Quisiera empezar dando una caracterización “de brocha gorda” de las propiedades de la mecánica cuántica, haciendo especial énfasis en las diferencias con la mecánica clásica. Las 3 propiedades siguientes me parecen las más importantes:

1. La mecánica cuántica es una teoría “todo o nada”. En otras palabras, al interactuar dos sistemas cuánticos, o no sucede nada, o por lo contrario sucede un efecto significativo.
2. Es una teoría *probabilista*. En combinación con 1., significa entre otros lo siguiente: lo que, en mecánica clásica, sería un efecto débil, en mecánica cuántica se vuelve un efecto fuerte pero poco probable. Por ejemplo, si una partícula pasa lejos de otra y la interacción es débil, en mecánica clásica la partícula sencillamente se desvía poco, mientras en la mecánica cuántica, se puede desviar arbitrariamente, pero con baja probabilidad.
3. Sin embargo, de cierta manera la teoría también es *determinista*: esto a primera vista parece contradecir 2., pero en realidad lo complementa. En mecánica cuántica existe un objeto, el *estado cuántico*, que permite realizar dos tareas: primero, predecir, para una medición arbitraria, las probabilidades que ésta arroje determinado resultado, y segundo, predecir el estado cuántico futuro del sistema para todos tiempos, siempre y cuando durante este tiempo no ocurra ninguna medición.

En términos formales, estas propiedades se entienden como sigue. El estado cuántico descrito en 3. es sencillamente el estado puro ψ del sistema, que es un vector *normalizado* de un espacio de Hilbert. Al conocer $\psi(0)$, puedo calcular $\psi(t)$ usando la ecuación de Schrödinger:

$$i\hbar\dot{\psi} = H\psi, \quad (1)$$

donde H es el Hamiltoniano del sistema, y se representa matemáticamente por un operador autoadjunto [19]. Integrando (1) obtenemos $\psi(t)$ para cada t conociendo $\psi(0)$. Por otro lado, como se dice en 3., se puede calcular la probabilidad de obtener un resultado dado midiendo una observable A arbitraria. Otra vez, A es un operador autoadjunto, con eigenvectores ϕ_n y eigenvalores a_n . Los a_n son los valores posibles que se pueden obtener de una medición de A , y si se mide A en cualquier estado ψ , la probabilidad que la medición arroje el resultado a_n está dada por la *regla de Born*

$$p_n = |\langle \psi, \phi_n \rangle|^2. \quad (2)$$

Aquí $\langle \psi, \phi_n \rangle$ se calcula mediante el *producto escalar*, que es parte de la definición de un espacio de Hilbert. Los p_n son, obvio está, forzosamente positivos. Que se suman a uno sigue de la *normalización* del estado ψ , que cumple la condición $\langle \psi, \psi \rangle = 1$. Son entonces legítimas probabilidades. Sería una buena tarea convencerse de esto, y también verificar que si el estado inicial está normalizado, sigue siéndolo para todos los tiempos.

*Electronic address: leyvraz@fis.unam.mx

Con ello nos convencimos de las observaciones 2. y 3. Queda 1. Pero esto es casi obvio: consideremos un estado inicial arbitrario $\psi(0) = \psi_0$, y definamos una base ortogonal del espacio de Hilbert que contiene ψ_0 como uno de sus elementos, digamos $\psi_0, \psi_1, \psi_2, \dots$. Entonces el estado $\psi(t)$ como función del tiempo se describe como

$$\psi(t) = \sum_{n=0} c_n(t) \psi_n. \quad (3)$$

Para t pequeños, resulta que la $c_0(t)$ es cercana a la unidad y los demás cercanos a 0. Ahora ¿cuál es el sentido de un estado como el descrito en (3)? Si preguntamos, por ejemplo, en cuál de los estados ψ_n se encuentra el estado $\psi(t)$, esto corresponde a la medición de una observable legítima. Pero aplicando la regla de Born (2), obtenemos que la probabilidad de encontrarse en el estado ψ_n está dada por $|c_n(t)|^2$, es decir, es baja si t es pequeño y $n \neq 0$, pero es distinta de 0. Sin embargo, el estado ψ_1 , por ejemplo, es radicalmente distinto del estado ψ_0 , ya que ambos son perpendiculares. Vemos entonces que existe una probabilidad que del estado ψ_0 pasemos al esencialmente distinto ψ_1 . Esto corresponde a la naturaleza “todo o nada” de la mecánica cuántica.

II. LA MEDICIÓN EN MECÁNICA CUÁNTICA

En la sección anterior, distinguimos entre la evolución de un estado puro en el tiempo mediante la ecuación de Schrödinger (1) por un lado, y la medición de cualquier cantidad cuando el sistema se encuentra en determinado estado por el otro. Fuimos bastante específicos, al decir que la ecuación de Schrödinger (1) sólo tiene validez mientras no se dé una medición del estado.

Esto ¡tiene mucho de raro! No debería aceptarse sin algo de molestia. Es profundamente distinto de lo que sucede en mecánica clásica: al fin y al cabo, Júpiter sigue las ecuaciones de Newton de la misma manera cuando lo mira un astrónomo, y cuando no. Las ecuaciones de Newton pretendían dar una descripción universal del mundo, y ésta incluía, obvio, la posibilidad de medir lo de que trata la teoría.

En mecánica cuántica no sucede así. No me atrevo, en modo alguno, a explicarles cuál es la solución “verdadera” del problema. Se observa que 5 físicos defienden, en promedio, entre 7 y 10 opiniones incompatibles al respecto, y la mayoría adicionalmente opina que el problema no tiene importancia. En lo último, por cierto, lo más probable es que tengan algo de razón: existe una variedad increíble de enfoques diferentes, pero ningún experimento los puede desempatar, por lo que el asunto se mantiene antes de todo en el ámbito de la teoría, o más aún de la filosofía.

Sin embargo sí se puede explicar de manera *parcial* por qué el problema tiene dificultad. Tiene que ver con la propiedad 1. de la mecánica cuántica, que es una teoría “todo o nada”. En efecto, resulta claro que, para averiguar algo de un sistema, tengo que entrar en contacto con él. Por ejemplo, en el caso sencillo de la observación visual, debo mandar luz sobre el objeto que quiero ver. Para un sistema clásico, siempre podemos usar lo bastante débil para que no afecte al objeto de nuestra observación. La vida diaria ofrece ejemplos con cierta similitud: si observamos alguien, y éste se da cuenta, es muy factible que se afecte su comportamiento, pero en principio siempre es posible observar a alguien sin que se dé cuenta. Es decir, siempre podemos observar el comportamiento de un sistema clásico sin que éste se vea afectado por dicha observación.

Debido a su naturaleza “todo o nada”, la situación está muy distinta en el caso de la mecánica cuántica. Por un lado, el aparato debe ser un objeto de tipo clásico: si no lo fuera, no podríamos decir que la misión de medir al sistema cuántico se hubiera cabalmente cumplido: en efecto, si el aparato es todavía un sistema cuántico, el problema de medir el resultado que dio el “aparato cuántico” es de la misma naturaleza que el problema de medir al sistema cuántico original. Si, por otro lado, el aparato de medición es clásico, entonces hemos terminado la tarea: el aparato arroja un resultado determinado; se puede ahora averiguar cuál es el resultado, ya que los sistemas clásicos sí se pueden observar sin perturbarlos.

Ciertamente hay algo un tanto perturbador en el hecho que no podemos dar cuenta del mundo entero en términos puramente cuánticos. Necesitamos un mundo clásico para observar al mundo cuántico, y necesitamos a la mecánica cuántica para explicar el mundo clásico en términos microscópicos. ¿Cómo resolvemos este malestar? No lo haré aquí: creo que es útil que el lector se quede con algunas preguntas sin respuesta: es un aliciente a pensar por cuenta propia. Mis opiniones al respecto no valen la pena describirse aquí.

La medición consiste, entonces, a conectar de tal manera un sistema cuántico a otro clásico, que el estado final del sistema clásico nos dé información acerca el estado cuántico original. En este sentido vemos una diferencia entre la medición y la evolución de un sistema cuántico que sigue la ecuación de Schrödinger (1). En el primer proceso intervienen por un lado un sistema cuántico, al que se mide, y por otro un sistema clásico que realiza dicha medición. En el segundo, sólo aparece un sistema cuántico, y no se puede registrar ningún resultado de medición.

Es ahora creíble, y un análisis más detallado lo confirma, que cualquier intento de acoplar un sistema de medición clásico—que debe tener la complejidad suficiente para registrar los resultados de la medición—a un sistema cuántico

afecta a éste último de manera incontrolable, de tal modo que al final sólo se puede hacer una predicción probabilista del valor que se medirá.

Cosas similares ocurren a veces en contextos más cotidianos: seguramente hemos visto, de niños, una bonita vajilla, y la hemos tocado para cerciorarnos de su realidad. Pero al tocarla, la vajilla se suele romper, por lo que se dice a los niños que “sólo miren con los ojos”. Los sistemas cuánticos son, en cierta medida, similares, con la salvedad que *no hay* manera segura de mirarlos. A pesar de todas las precauciones que se toman, si queremos averiguar algo del sistema cuántico, debemos romperlo [20].

Esto es una parte importante de la dificultad que hay en medir sistemas cuánticos. Resulta, sin embargo, que no es para nada lo único. Si tenemos un sistema que consiste de 2 partes, entonces al medir una parte, puedo afectar de alguna manera a la otra parte de manera instantánea. Obviamente, semejante “influencia a distancia” no es algo que puede explicarse por interacciones físicas del tipo habitual. La existencia y la naturaleza de estas influencias son el objeto del resto de este artículo.

III. DEFINICIÓN DE LOCALIDAD

El tema central de las siguientes secciones será el de mediciones distantes y simultáneas sobre partes distintas de un mismo sistema. Definamos primero el sentido que daremos a estas palabras.

Lo más importante es tal vez la palabra “simultánea”. Tenemos un sistema que consta de 2 partes, la una se encuentra en la región 1 y la otra en la región 2. Supondremos que el proceso de medición toma un tiempo corto, y que ambas regiones están bien separadas. Es entonces posible suponer que dichas mediciones estén tan lejanas en el espacio, y tan cercanas en el tiempo, que un rayo de luz emitido en 1 al tiempo de la primera medición no alcance 2 al tiempo de la segunda. En este caso, decimos que las mediciones son *simultáneas*. En este caso sabemos, en base a la relatividad especial de Einstein, que nada de lo que sucede en la medición en 1 puede afectar la medición en 2, y viceversa.

En realidad, consideraremos una situación más restrictiva aún. Como vimos, la elección de la cantidad que se va a medir afecta fuertemente al desarrollo del estado cuántico durante el proceso de medición. Por lo tanto, no sólo vamos a suponer que las mediciones mismas tengan lugar “simultáneamente”, sino que las decisiones acerca de las cantidades que se medirán se toman también simultáneamente. En otras palabras, todo el proceso de determinar la cantidad que se va a medir, seguido por la medición misma, tiene lugar en cada una de las regiones 1 y 2, y dura un tiempo menor del que un rayo de luz requiere para ir de 1 a 2.

Finalmente debemos decir a qué nos referimos al decir que medimos 2 partes “de un mismo sistema”. La situación típica que se considera es una en que se producen 2 o más partículas en algún estado que incluye correlaciones entre éstas. Las diferentes partículas se envían luego a lugares distantes, las regiones 1 y 2 arriba mencionadas, dónde se realizan varias mediciones sobre cada partícula. Las mediciones se hacen simultáneamente en el sentido que hemos escrito arriba: las mediciones se eligen independientemente y al azar en cada región. Esto significa que la partícula en la región 2 no puede “saber” qué medición se realizó sobre la partícula en la región 1, y viceversa.

Bajo estas circunstancias, hacemos la pregunta siguiente: si en la región 1 realizo al azar 2 mediciones posibles distintas, digamos A_1 y B_1 (en lo siguiente, las letras indicarán siempre un determinado tipo de medición, mientras el índice nos dirá en qué región se realiza). Por otro lado supongamos que en la región 2 hacemos siempre la misma medición A_2 . ¿Será posible que el comportamiento estadístico de los resultados de las mediciones de A_2 sean distintos según se esté midiendo A_1 o B_1 en la región 1?

Creo que es aparente que no puede ser así: en efecto, de ser el caso, sería posible mandar una señal de la región 1 a la región 2 eligiendo medir B_1 en vez de A_1 , lo que contradice a la teoría especial de la relatividad. Veamos como puede hacerse: supongamos, por ejemplo, que todas las mediciones arriba mencionadas siempre dan 0 o 1 como resultado. Supongamos adicionalmente que si mido A_1 en la región 1, la cantidad A_2 da 0 o 1 con probabilidad $1/2$, mientras que si mido B_1 , la probabilidad de medir A_2 igual a 1 crece a 0,6. Esto es un ejemplo típico del efecto que describí en términos generales arriba.

¿Podemos en este caso realizar una señalización instantánea de la región 1 a la 2? A primera vista no parece: si en la región 2 mido $A_2 = 1$, esto no me garantiza de modo alguno que se haya medido la B_1 en la región 1: la diferencia entre $1/2$ y $0,6$ no permite este tipo de conclusiones. Pero podemos tomar un número mayor de sistemas: si tomamos, digamos, 1000 sistemas del mismo tipo, y sometemos todas las partículas de los 1000 sistemas que se encuentran en 1 a la medición B_1 , entonces esperamos que en la región 2 tendremos aproximadamente 600 valores de A_2 igual a 1, mientras con una medición en la región 1 de A_1 , sólo tendríamos 500 valores de A_2 iguales a 1. Semejante diferencia sí permite al observador en la región 2 determinar con un muy alto grado de probabilidad si se midió A_1 o B_1 en la región 1.

Generalizando las observaciones que acabo de hacer, debería quedar claro que, en caso de existir diferencias estadísticas en los resultados de las mediciones en la región 2, según las cantidades que se hayan medido en la región 1,

sería posible usarlas para mandar información sobre distancias arbitrarias, entrando así en conflicto con la relatividad.

Es por lo tanto un gran alivio enterarse que la mecánica cuántica en este sentido es local, véase [1, 2] por una prueba sencilla, pero rigurosa. Formulemos esto explícitamente:

Localidad Estadística: Si se repiten muchos experimentos en los que se realizan mediciones distantes y simultáneas sobre partes del mismo sistema, las mediciones de una cantidad en una región tienen el mismo comportamiento estadístico independientemente de las mediciones que se hacen en otra región.

Esto entonces nos dice que la mecánica cuántica no es tan extraña como a veces se afirma. El concepto más elemental de localidad está firmemente anclado en su formalismo. Adicionalmente, las pruebas de la localidad, por ejemplo en [1] no son de matemáticas muy avanzadas con recónditas sutilezas, sino que son más bien sencillas y convincentes. ¿Por qué, entonces, seguimos escuchando de la supuesta “no-localidad” de la mecánica cuántica? En lo siguiente no seguiremos el camino de Bell y otros, sino nos atenderemos a un enfoque descrito en particular por Peres en [3].

Se trata de una posible extensión de la *Localidad Estadística* a algo más fuerte que, a falta de mejor palabra, llamaré la *Localidad Hipotética*, aunque tal vez éste no sea un nombre muy atractivo. Notemos que ésta no es terminología estándar, pero ello no es muy grave, ya que la terminología generalmente aceptada en este tema es relativamente compleja. Veamos de qué se trata.

Un limitante aparente de la Localidad Estadística es que sólo se refiere a una multiplicidad de experimentos independientes, y al comportamiento estadístico de las mediciones en la región 2, que no se ve afectado por las mediciones en la región 1. Podríamos desear más: por ejemplo, que el *verdadero estado de cosas* en la región 2 no se vea afectado por un cambio en la cantidad que se midió en la región 1. Formularé una versión de esta creencia que parecerá a primera vista algo extraña, pero que seguramente capta cierta parte de nuestra intuición sobre las consecuencias de la imposibilidad de influencia a distancia:

Localidad Hipotética: Si se hace un solo experimento en el que se realizan mediciones distantes y simultáneas sobre partes del mismo sistema, la medición de una cantidad en una región *hubiera tenido* el mismo resultado que el que efectivamente tuvo, si en las otras regiones *se hubieran realizado* mediciones de cantidades distintas.

La primera reacción bien puede ser, que quien discute semejante concepto se la fumó verde. ¿Cómo podemos saber lo que hubiera pasado en un caso que no se realizó? Y por otro lado ¿por qué tendría esto la menor importancia?

A estas preguntas hay varias respuestas. Primero adelantamos el resultado que da a este concepto cierta relevancia: mientras es a todas luces obvio que nunca podremos verificar el postulado de la Localidad Hipotética (por la evidente imposibilidad de realizar un experimento de dos maneras contradictorias) resulta que sí es posible que un experimento muestre resultados en conflicto con este postulado. En otras palabras, es imposible confirmarlo, pero no es imposible refutarlo. Y en efecto la mecánica cuántica, como veremos en la Sección IV, predice resultados que contradicen a este principio. Éstos pueden también verificarse experimentalmente

Otro motivo de interesarse a esta propiedad, a primera vista algo extraña, es la existencia de una gran clase de teorías para las que se cumple. Por ejemplo, en mecánica clásica, y generalmente hablando en todas las teorías “clásicas”, es verdad que, si cambiáramos hipotéticamente algo en un lugar, nada cambiaría a distancia suficiente: se pueden hacer predicciones teóricas sobre realidades alternas, y todas las teorías clásicas, en las que la observación no juega ningún papel especial, predicen que la Localidad Hipotética, no es una aberración, sino más bien una trivialidad. En efecto, imaginémosnos que las 2 partes del sistema están descritas con una cantidad de poleas, correas, campos electromagnéticos y otros más objetos de la más variada índole, descritos por algún tipo de teoría clásica (mecánica para las poleas, electrodinámica al estilo de Maxwell para los campos electromagnéticos, y la teoría favorita del lector, mientras ésta sea de tipo clásico y local, para lo demás). Entonces la teoría puede usarse para contestar la pregunta de lo que *hubiera* pasado en la región 1, si en la región 2 se *hubieran* realizado mediciones distintas. Como las teorías, por hipótesis, son locales, en el sentido de cumplir con la relatividad de Einstein, es obvio que si cambiamos algo en las mediciones que tienen lugar en la región 2, y si éstas son simultáneas con las mediciones en la región 1, la teoría predecirá que no habrá efecto.

Tampoco podemos pasar por alto el atractivo *intuitivo* que tiene este postulado. Si no puede haber influencia instantánea de un lugar a otro, entonces ¿acaso no está *obvio* que cualquier cosa que suceda en la región 1 no puede afectar a lo que sucede en la región 2, sin importar cuál alternativa realmente haya sucedido? Parte de esta convicción viene en parte de la íntima conexión entre alternativas que no ocurrieron y nuestras ideas de causalidad: al decir que un evento *A* causó otro evento *B*, muchas veces pensamos que, de haberse evitado *A*, posiblemente *B* no hubiera sucedido. Y viceversa, si estamos convencidos que el evento *A* no tuvo ninguna parte en que *B* aconteciera, se traduce esta creencia en la firme convicción que *B* hubiese ocurrido de la misma manera, también si *A* no hubiera ocurrido. Como ejemplo de lo último podemos pensar en el resultado de un experimento físico, digamos un gran experimento del CERN, y como evento *A* la fase de la Luna durante el experimento. Decir que la fase de la Luna *no influyó* en el experimento, parece muy equivalente a la afirmación que, si la Luna hubiera tenido una fase distinta, el experimento hubiera salido igual.

Pero entonces, al evidenciar que la mecánica cuántica no cumple con esta propiedad, mostramos algo que sí tiene cierto interés. En particular, si falla la Localidad Hipotética en la mecánica cuántica, resulta que los resultados

de la mecánica cuántica no se pueden explicar de ninguna manera por cualquier teoría de naturaleza clásica, por compleja que ésta pudiera resultar. En otras palabras, las teorías que no cumplen con la Localidad Hipotética, son de naturaleza esencialmente diferente de todas las teorías de tipo clásico, o usando una expresión común en el campo, son muy distintas de todas las teorías realistas.

Pero prometí enseñar violaciones de la Localidad Hipotética en la mecánica cuántica: es tiempo que cumpla.

IV. LA LOCALIDAD HIPOTÉTICA NO SE CUMPLE EN MECÁNICA CUÁNTICA

Daremos 2 ejemplos de sistemas en los que se viola la Localidad Hipotética. Aquí no entraremos en los detalles que prueban que, en efecto, la mecánica cuántica permite crear sistemas con las desconcertantes propiedades que a continuación describo. Sin embargo, esta prueba se realiza en unos apéndices, que requieren un poco más de conocimientos del formalismo de la mecánica cuántica.

El primer sistema es una versión de un experimento propuesto por Lucien Hardy [4] y funciona como sigue (seguimos aquí la admirable presentación pedagógica de Mermin [5]). Como siempre, tenemos dos regiones 1 y 2 donde se realizan dos tipos de mediciones, A y B , cuyos resultados pueden ser 0 o 1. El experimento se hace muchas veces, y en cada realización se escoge en cada región si se hace la medición A o la B al azar y de manera independiente. Después se recogen los datos y se cotejan los resultados: ¿qué resultados obtenemos si medimos A_1 y B_2 juntos? ¿ B_1 y A_2 ? ¿Cuál es la lista de todas estas correlaciones?

En el apéndice A mostramos que, usando propiedades específicas de la mecánica cuántica, es posible arreglar las cosas de tal manera que se cumplan las *reglas* siguientes

AA Si ambas mediciones son A , entonces a veces ambas dan el valor 1. Para ser específicos, se puede llegar a que esto suceda 9% de las veces.

AB Si una de las mediciones es A y la otra B , *nunca* ambas mediciones son iguales a 1.

BB Si ambas mediciones son B , *nunca* ambas mediciones son iguales a 0.

Estas 3 condiciones obviamente no encierran ninguna contradicción. Es más, no existe ninguna dificultad en realizarlas si el resultado en la región 1 puede influir en él de la región 2. Pero esto es precisamente lo que no queremos. Que se puede diseñar un sistema físico cumpliendo estas características se muestra en el Apéndice A.

Veamos ahora si estas condiciones, que a primera vista no parecen especialmente maléficas, son compatibles con la Localidad Hipotética. Consideremos uno de los casos con ambos valores de A iguales a 1: por la condición **AA** sabemos que éstos ocurren. Ahora aplicamos la Localidad Hipotética a este experimento en particular. En este caso lo que realmente medimos fueron A_1 y A_2 , y los resultados reales fueron 1 en ambos casos. ¿Qué hubiera pasado si hubiera medido B_2 en vez de A_2 ? Claramente, por la Localidad Hipotética, el valor de A_1 que se hubiera medido hubiera sido el mismo que el que se midió en el mundo real, es decir 1. Pero en esta medición hipotética, las reglas arriba indicadas se deben seguir cumpliendo. Ya que el valor de A_1 hubiera sido 1, por la regla **AB** es necesario que la B_2 hubiera dado un resultado de 0. Similarmente, ya que existe una simetría perfecta entre las regiones 1 y 2, nos podemos convencer que si hubiéramos medido B_1 en vez de A_1 , el B_1 hubiera arrojado un valor de 0.

Ahora, si lees esto cuidadosamente, ves que estamos ya peligrosamente cerca de una contradicción. En efecto hemos mostrado que, al cambiar A por B en cualquier lado, obtenemos con seguridad que la medición de B dará 0. Esto sugiere que, si hubiéramos medido la B en *ambos* lados, hubiéramos obtenido 0 de ambos lados, en manifiesta contradicción con la regla **BB**. ¿Podemos deducir esto de la Localidad Hipotética? Sí, pero vamos a requerir un nivel de abstracción adicional. Si empezamos con el experimento real, con una medición de A_1 y A_2 dando el valor 1 para ambos, entonces podemos primero cambiar A_1 a B_1 . Por el argumento arriba indicado, B_1 tiene que dar el valor 0. Ahora volvemos a cambiar A_2 por B_2 : debido a la regla **BB**, B_2 tiene que ser 1. Pero ahora puedo cambiar otra vez B_1 a A_1 . Que A_1 da el valor 1 es innegable, de tal manera que nos encontramos con una situación en la que B_2 da el valor 1 y A_1 también, en contradicción con la regla **AB**.

Dicho de otro modo, no existe ninguna manera de dar valores simultáneas a A_1 , A_2 , B_1 , y B_2 que cumplan con las tres reglas: en efecto, A_1 y A_2 pueden tener ambos el valor 1. Para estos casos la regla **AB** impone el valor 0 tanto a B_1 como a B_2 . Pero que ambos B_1 y B_2 tengan el valor 0 contradice a **BB**.

Este ejemplo no es único, sino que hay una gran variedad de fenómenos similares. Veamos uno más. Es posible mandar 3 partículas correlacionadas en 3 direcciones distintas, y una vez que las 3 partículas estén lo bastante lejos, realizar al azar una de 2 mediciones, que para variar llamamos X e Y , como siempre con un índice para indicar el lugar donde se realiza la medición. El resultado de la medición, como siempre, es un sí o un no, o en otras palabras un 1 o un 0. Ahora, después de llevar a cabo el experimento muchas veces, observamos las regularidades siguientes:

XXY Si se hacen dos mediciones de X y una de Y —en otras palabras, si se mide (X_1, X_2, Y_3) o (X_1, Y_2, X_3) o (Y_1, X_2, X_3) —el número total de respuestas afirmativas es impar.

YYY Si todas las mediciones son de Y , el número total de respuestas afirmativas es par.

Hagamos ahora la pregunta: ¿será posible determinar respuestas $X_{1,2,3}$ e $Y_{1,2,3}$ de manera que se cumplan las dos reglas arriba descritas? Mostremos que no: en las tres series (X_1, X_2, Y_3) , (X_1, Y_2, X_3) y (Y_1, X_2, X_3) , hay en total un número impar de respuestas afirmativas: en efecto, hay un número impar de éstas en cada serie, por la regla **XXY**, y hay 3 series: y ¡sumar 3 número impares da un número impar! Pero en esta serie la totalidad de los X tiene un número par de respuestas afirmativas, ya que cada X aparece 2 veces. Vemos entonces que las mediciones restantes tienen un número impar de respuestas afirmativas. Pero estas mediciones son las (Y_1, Y_2, Y_3) , y éstas tienen un número par de respuestas por la regla **YYY**. Esta contradicción indica que no podemos determinar una serie de respuestas $X_{1,2,3}$ e $Y_{1,2,3}$ que cumplan con las dos regularidades. Si esta prueba no te convence, puedes buscar todas las posibilidades; hay un total de $2^6 = 64$, y puedes convencerte en menos de 20 minutos que ninguna de las 64 posibilidades cumple con ambas reglas.

Ahora mostramos que, si suponemos la Localidad Hipotética, es necesario poder determinar los tres valores de X y de Y , e contradicción con lo que acabamos de mostrar. Para ello, consideramos una medición de X_1, Y_2 y Y_3 . Esto es un caso en que las reglas no nos dicen si las respuestas afirmativas están en número par o impar. Sin embargo, ya que se realizó el experimento en esta configuración, sabemos cuáles son los valores de X_1, Y_2 y Y_3 . Ahora preguntamos ¿qué hubiera pasado si hubiéramos medido Y_1 en vez de X_1 ? Estaríamos en un caso donde vale la regla **YYY**, lo que permite determinar con seguridad el valor de Y_1 que se hubiera medido: es el que da un número par de respuestas afirmativas al combinarse con Y_2 y Y_3 , que ya están conocidos. Similarmente, si hubiéramos medido X_2 en lugar de Y_2 , estaríamos en un caso donde se aplica la regla **XXY**, lo que permite determinar a su vez el valor que hubiera tenido X_2 . Y finalmente se logra lo mismo para X_3 si hubiera medido éste en vez de Y_3 . En otras palabras, hubiéramos podido determinar sin ambigüedad a los 6 valores $X_{1,2,3}$ e $Y_{1,2,3}$, lo cual, como acabamos de ver, es imposible.

Otra vez podemos preguntarnos si verdaderamente se puede realizar un sistema como el que acabo de describir. En el Apéndice B muestro, usando algo del formalismo de la mecánica cuántica, que ésta sí predice tales comportamientos.

V. SEPARABILIDAD Y DESIGUALDADES DE BELL

El enfoque que presenté hasta ahora es el que me parece más claro, pero no es el más tradicional. Para ser completo, quiero esbozar rápidamente el enfoque que se debe originalmente a Bell [6, 7] y que Mermin ha explicado de manera muy sencilla en [8–10]. Se centra en un concepto afín a, pero algo distinto de, la Localidad Hipotética, que se llama la separabilidad.

Consideremos varias mediciones posibles en las regiones 1 y 2, digamos A y B . Sin pérdida de generalidad podemos suponer que los únicos resultados posibles a estas mediciones son respuestas sí o no, o de manera equivalente, 1 o 0. La correlación $p_{++}(A_1, A_2)$, por ejemplo, es entonces la probabilidad de obtener dos respuestas afirmativas a las preguntas A_1 y A_2 respectivamente.

Si las mediciones se hacen sobre sistemas que no tienen ninguna relación entre sí, entonces las mediciones son *independientes*, es decir

$$p_{++}(A_1, A_2) = p_+(A_1)p_+(A_2) \quad (4)$$

donde $p_+(A_1)$ es la probabilidad de obtener una respuesta positiva a A_1 independientemente de lo que pasa en la región 2, y similarmente para $p_+(A_2)$. Pero obviamente es posible, sin problema alguno, tener correlaciones entre las 2 mediciones. Las 2 partículas, al llegar a sus lugares de medición, pueden traer consigo información que viene de su pasado común, expresada en el valor de algún parámetro λ . Podemos entonces obtener correlaciones distintas de (4) si suponemos que los valores medidos de A_1 y A_2 dependen de λ y estén dados por $A_1(\lambda)$ y $A_2(\lambda)$. Por ejemplo, en el caso extremo en que, para cada valor de λ , $A_1(\lambda) = A_2(\lambda)$, entonces tenemos

$$p_{++}(A_1, A_2) = 1. \quad (5)$$

Sin embargo, los valores individuales de A_1 y A_2 fluctúan en base a las variaciones de λ , de modo que $p_+(A_1) = p_+(A_2) \neq 1$.

Estas consideraciones sugieren la representación siguiente para las correlaciones de este tipo: una vez dados todos los elementos que pudieran contribuir de manera común a las mediciones en ambas regiones, y una vez resumidos todos estos elementos en una variable (posiblemente vectorial) λ , están determinados en función de λ los resultados de las varias mediciones posibles en la región 1, como $A_1(\lambda)$, $B_1(\lambda)$ etc. así como también los resultados de las mediciones en la región 2. Las correlaciones más allá de la independencia se deben entonces, se supone, a la existencia de un

pasado común que influye sobre ambos resultados. Este parámetro que describe el pasado común toma el valor λ con una probabilidad $\rho(\lambda)d\lambda$. Por lo tanto, decimos que

$$p_{++}(A_1, A_2) = \int d\lambda \rho(\lambda) A_1(\lambda) A_2(\lambda) \quad (6)$$

Por otro lado, ya que las mediciones son *simultáneas*, la medición de A_1 en la región 1 no puede depender en modo alguno de lo que sucede en la región 2. Por ello, si mido B_2 en lugar de A_2 en la región 2, obtengo

$$p_{++}(A_1, B_2) = \int d\lambda \rho(\lambda) A_1(\lambda) B_2(\lambda) \quad (7)$$

donde $A_1(\lambda)$ es *la misma función* que la que usamos en (6). Ésta es la propiedad fundamental de localidad, o separabilidad, que permite obtener los resultados que abajo se detallan.

Esta serie de hipótesis se conocen como la hipótesis de separabilidad, localidad de Bell o localidad de Einstein. Un poco de reflexión mostrará que está muy relacionada con la de la Localidad Hipotética: los $A(\lambda)$ que no se midieron no son otra cosa que resultados de experimentos que no se llevaron a cabo. Sin embargo, es un planteamiento matemáticamente algo más limpio y se puede usar para obtener resultados matemáticos rigurosos cerca de las probabilidades $p_{++}(A_1, A_2)$.

Estos resultados toman a menudo la forma de desigualdades, conocidas de manera genérica como desigualdades de Bell. Una derivación sencilla descansa sobre el hecho elemental que, si x e y toman los valores ± 1 , entonces

$$0 \leq |x + y| + |x - y| \leq 2. \quad (8)$$

Para probarlo ¡sustituye las 4 posibilidades! Con un poco más de trabajo se puede extender a todos los valores $0 \leq x, y \leq 1$. Por ende vemos que para cada λ necesariamente tenemos

$$\begin{aligned} 0 &\leq |A_1(\lambda)A_2(\lambda) + A_1(\lambda)B_2(\lambda) + B_1(\lambda)A_2(\lambda) - B_1(\lambda)B_2(\lambda)| \\ &\leq |A_1(\lambda)[A_2(\lambda) + B_2(\lambda)]| + |B_1(\lambda)[A_2(\lambda) - B_2(\lambda)]| \\ &\leq |A_2(\lambda) + B_2(\lambda)| + |A_2(\lambda) - B_2(\lambda)| \leq 2 \end{aligned} \quad (9)$$

Pero si la desigualdad (9) es válida para cada λ , entonces seguirá siendo válida si realizo un *promedio* sobre todos los λ , como se plantea en (6,7). De esto se puede concluir que, para cualesquieras mediciones A_1, A_2, B_1 y B_2 tenemos

$$p_{++}(A_1, A_2) + p_{++}(A_1, B_2) + p_{++}(B_1, A_2) - p_{++}(B_1, B_2) \leq 2. \quad (10)$$

Ésta es una de las múltiples variantes de las desigualdades de Bell [6, 7], también conocida como desigualdad de CHSH (Clauser–Horne–Shimony–Holt) [11, 12]. En muchos casos la mecánica cuántica prevé resultados que violan estas desigualdades. Como puede verse en [11, 12], es posible llevar a cabo experimentos que permitan decidir si estas predicciones de la mecánica cuántica son correctas o no. Hasta ahora no ha habido experimentos convincentes confirmando las desigualdades, y muchos por otro lado que sí han confirmado el resultado de la mecánica cuántica. Por lo tanto debemos considerar que la mecánica cuántica realmente no cumple con los postulados que se requieren para derivar (10). Los trabajos experimentales clásicos son los de Alain Aspect [13, 14]. Un experimento llamativo por haberse realizado sobre distancias especialmente largas se hizo en [15].

VI. CONCLUSIONES

Y ahora ¿cuál es la palabra final? ¿Qué hemos de creer? ¿Realmente sigue de todo aquello que en el mundo “todo influye en todo”? ¿Realmente habrá influencias a distancias arbitrarias que conectan todo? Hay muchas respuestas, y para saber más, refiero el lector interesado a varios artículos que van más lejos, véanse por ejemplo [16–18], así como las referencias que allá se encuentran.

No quiero concluir de manera dogmática. Existen varias maneras de acostumbrarse a la idea que la Localidad Hipotética es inválida. Una es sencillamente recordar siempre lo dicho por Asher Peres [3]: “Unperformed experiments have no results”. Las contradicciones, la posible influencia no-local a distancia arbitraria, sólo se encuentran al hacer preguntas que, por su naturaleza misma, no pueden decidirse por la observación, al referirse a experimentos que no se llevaron a cabo. Si recordamos que es parte esencial de la física el apoyarse en resultados experimentales, podemos tal vez concluir que los fenómenos que hemos discutido aquí no tienen mayor importancia en nuestro quehacer como físicos.

Por otro lado, el propósito de la física seguramente no debe excluir preguntas fundamentales sobre la naturaleza de la realidad que nos rodea. La pregunta de saber si realmente existen, o no, influencias no-locales en el mundo real es de gran importancia, y tal vez aún sin tener elementos experimentales para dar una respuesta, alguna especulación sea permitida. Lo que vimos es lo siguiente: si en una región 1, hubiéramos hecho un experimento diferente al que realmente hicimos, sabemos que los resultados en la lejana región 2 sí hubieran resultado afectados. Pero ¿de qué manera cambiarían? Ningún experimento nos lo puede decir. Sin embargo, algo sí sabemos, por la Localidad Estadística: todas las propiedades estadísticas de las mediciones realizadas en la región 2 hubieran permanecido idénticas. En otras palabras, si las mediciones en la región 2 son, por ejemplo, una serie aleatoria de 0 y 1 con una probabilidad de 1/2 para ambos, se mantendrían estas características en cada variación hipotética del experimento. En otras palabras, esta variación no puede crear ningún efecto, y es inapropiado preguntar qué hubiera sido la “causa” de la variación. Más bien ésta ocurre, o mejor dicho ocurriría, fuera de cualquier contexto causal. Sólo se trata de sustituir una instancia de un proceso azaroso por otra, prácticamente indistinguible. Bajo estas circunstancias, podemos, tal vez, aceptar que al malestar que nos puedan causar estas “no-localidades” no se le tiene que dar excesiva importancia.

Por otro lado, si insistimos en imponer una estructura “realista” a la mecánica cuántica, entonces no puede caber duda que los elementos que determinan esta realidad tienen forzosamente que tener conexiones instantáneas, no-locales y cuya intensidad no disminuye con la distancia. Estos atributos pueden parecer algo perturbadores, pero al aplicarse sólo a objetos absolutamente inobservables, no causan mayores trastornos en la física propiamente dicha [21].

Apéndice A: La paradoja de Hardy

Aquí mostraremos como se puede realizar la paradoja de Hardy. Tenemos 2 partículas, y éstas se deben poder medir para dar resultados de tipo 1 o 0. 2 partículas que viven en un espacio de Hilbert de 2 dimensiones son entonces apropiadas, es decir, el espacio de Hilbert del sistema entero es el producto tensorial de 2 espacios de Hilbert bidimensionales, es decir, tenemos un espacio de 4 dimensiones.

Cada uno de los operadores $A_{1,2}$ y $B_{1,2}$ genera su propia base ortogonal de eigenvectores. Definamos los ϕ como eigenvectores de los A y los ψ como los de los B .

$$A_1\phi_+^{(1)} = \phi_+^{(1)} \quad (\text{A1a})$$

$$A_2\phi_+^{(2)} = \phi_+^{(2)} \quad (\text{A1b})$$

$$B_1\psi_+^{(1)} = \psi_+^{(1)} \quad (\text{A1c})$$

$$B_2\psi_+^{(2)} = \psi_+^{(2)} \quad (\text{A1d})$$

mientras los eigenvectores correspondientes con índice $-$ corresponden al eigenvalor 0.

Ahora debo elegir un estado Ψ del sistema entero que cumpla con las varias reglas descritas en el texto. Para cumplir primero con la regla **BB**, el estado debe ser de la forma

$$\Psi = \alpha\psi_+^{(1)} \otimes \psi_-^{(2)} + \beta\psi_-^{(1)} \otimes \psi_+^{(2)} + \gamma\psi_+^{(1)} \otimes \psi_+^{(2)}. \quad (\text{A2})$$

Las condiciones **AB** se se expresan como

$$(\psi_+^{(1)} \otimes \phi_+^{(2)}, \Psi) = 0 \quad (\text{A3a})$$

$$(\phi_+^{(1)} \otimes \psi_+^{(2)}, \Psi) = 0 \quad (\text{A3b})$$

o equivalentemente

$$(\phi_+^{(2)}, \alpha\psi_-^{(2)} + \gamma\psi_+^{(2)}) = 0, \quad (\text{A3c})$$

$$(\phi_+^{(1)}, \beta\psi_-^{(1)} + \gamma\psi_+^{(1)}) = 0. \quad (\text{A3d})$$

Ahora es claro que, para cada valor de α , β y γ , estas condiciones se pueden cumplir al elegir

$$\phi_+^{(1)} = \frac{1}{\sqrt{\beta^2 + \gamma^2}} \left[-\beta\psi_-^{(1)} + \gamma\psi_+^{(1)} \right], \quad (\text{A4a})$$

$$\phi_+^{(2)} = \frac{1}{\sqrt{\alpha^2 + \gamma^2}} \left[-\alpha\psi_-^{(2)} + \gamma\psi_+^{(2)} \right]. \quad (\text{A4b})$$

En otras palabras, para cada estado de la forma (A2) se puede elegir unos A_1 , B_1 , A_2 y B_2 que cumplan las reglas de la paradoja de Hardy.

Finalmente evaluemos la probabilidad de tener que tanto A_1 como A_2 den respuesta afirmativa. Ésta se da por

$$p_{AA} = \left| (\Psi, \phi_+^{(1)} \otimes \phi_+^{(2)}) \right|^2. \quad (\text{A5})$$

Un cálculo elemental muestra que

$$p_{AA} = \frac{|\alpha\beta\gamma|^2}{(|\alpha|^2 + |\gamma|^2)(|\beta|^2 + |\gamma|^2)} \quad (\text{A6})$$

con la condición adicional debida a la normalización de Ψ :

$$|\alpha|^2 + |\beta|^2 + |\gamma|^2 = 1. \quad (\text{A7})$$

Un cálculo algo engorroso, pero que no tiene mayor dificultad muestra que el valor máximo de p_{AA} para valores de α , β y γ cumpliendo (A7) está dado por

$$p_{AA} \leq \left(\frac{\sqrt{5} - 1}{2} \right)^5 \approx 0,09017 \quad (\text{A8})$$

lo cual justifica lo dicho en el texto, que se puede lograr que A_1 y A_2 den el resultado (1, 1) en aproximadamente 9% de los casos.

Apéndice B: El estado de Greene–Hornberger–Zeilinger (GHZ)

El ejemplo que trataremos aquí se obtiene usando las propiedades de las matrices de Pauli, que se definen como sigue:

$$\sigma_x = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}, \quad \sigma_y = \begin{pmatrix} 0 & -i \\ i & 0 \end{pmatrix}, \quad \sigma_z = \begin{pmatrix} 1 & 0 \\ 0 & -1 \end{pmatrix}. \quad (\text{B1})$$

Las 2 propiedades esenciales que usaremos en lo siguiente son las siguiente: primero, 2 matrices de Pauli distintas siempre *anticommutan*; segundo, el cuadrado de cualquiera de las matrices de Pauli es siempre la identidad. En fórmulas

$$\sigma_i \sigma_j = -\sigma_j \sigma_i \quad (i \neq j) \quad (\text{B2a})$$

$$\sigma_i^2 = \mathbb{1} \quad (\text{B2b})$$

También conviene notar que las matrices de Pauli son todas matrices hermitianas, de manera que se pueden usar para definir observables.

Ahora definimos los 4 operadores siguientes en el triple producto tensorial de un espacio de Hilbert bidimensional consigo mismo.

$$T_1 = \sigma_y \otimes \sigma_x \otimes \sigma_x, \quad (\text{B3a})$$

$$T_2 = \sigma_x \otimes \sigma_y \otimes \sigma_x, \quad (\text{B3b})$$

$$T_3 = \sigma_x \otimes \sigma_x \otimes \sigma_y, \quad (\text{B3c})$$

$$T_0 = \sigma_y \otimes \sigma_y \otimes \sigma_y. \quad (\text{B3d})$$

De (B2a) siguen inmediatamente las 2 propiedades siguientes de los T_k

1. Todos los T_k conmutan entre sí: esto sigue del hecho que, al conmutar 2 T 's, siempre aparecen dos anticonmutaciones que llevan a que los operadores conmuten.
2. T_0 tiene la representación siguiente

$$T_0 = -T_1 T_2 T_3 \quad (\text{B4})$$

Ahora, si 4 operadores conmutan, es bien conocido que se puede escoger un eigenvector común. Por simetría es razonable suponer que para este eigenvector los 3 operadores T_1 , T_2 y T_3 tendrán el mismo eigenvalor. Adicionalmente, ya que $T_k^2 = \mathbb{1}$ para cada k , vemos que los eigenvalores de los T_k sólo pueden ser ± 1 .

Si escogemos Ψ tal que

$$T_k \Psi = \Psi \quad (1 \leq k \leq 3) \quad (\text{B5})$$

entonces está claro que

$$T_0 \Psi = -\Psi. \quad (\text{B6})$$

Si el sistema está en el estado Ψ , entonces resulta que cualquier medición de los operadores T_1 , T_2 y T_3 dan el valor, digamos, 1, mientras cualquier medición de T_0 da -1 . Si entonces llamamos una medición de $\sigma_x \otimes \mathbb{1} \otimes \mathbb{1}$ una medición de X_1 , una de $\mathbb{1} \otimes \sigma_y \otimes \mathbb{1}$ una medición de Y_2 y así sucesivamente, entonces vemos que se cumplen las reglas que observamos en el sistema GHZ, tal y como se describe en el texto.

Como observación final hacemos notar que el origen de problema radica precisamente en la posibilidad de tener anticonmutación entre matrices, y al mismo tiempo de compensarlas en productos de varios términos. Si se cumpliera lo requerido por la Localidad Hipotética, o por la separabilidad, deberíamos poder asignar valores a cada uno de los σ_x y σ_y . Pero que esto es imposible, ya que, si los σ 's tuvieran valores, tendrían que ser ± 1 y seguiría de inmediato

$$T_0 = T_1 T_2 T_3 \quad (\text{B7})$$

en contradicción absoluta con (B4)

-
- [1] G.C. Ghirardi, A. Rimini y T. Weber, A general argument against superluminal transmission through the quantum mechanical measurement process, *Lettere Al Nuovo Cimento* (1971–1985), **27** (10), 293–298 (1980).
- [2] G.C. Ghirardi y T. Weber, Quantum mechanics and faster-than-light communication: methodological considerations, *Il Nuovo Cimento B* (1971–1996) **78** (1) 9–20 (1983)
- [3] A. Peres, Unperformed experiments have no results, *Am. J. Phys.* **46** (7) 745–747 (1978)
- [4] L. Hardy, Quantum mechanics, local realistic theories, and Lorentz-invariant realistic theories, *Phys. Rev. Lett.* **68**, 2981–2984 (1992)
- [5] N.D. Mermin, Quantum mysteries refined, *Am. J. Phys.* **62** (10) 880–887 (1994)
- [6] J.S. Bell, On the Einstein-Podolsky-Rosen paradox, *Physics* **1** (3) 195–200 (1964)
- [7] J.S. Bell, Bertlmann's socks and the nature of reality, *J. Physique* **C2** 41–62 (1981).
- [8] N.D. Mermin, Bringing home the atomic world: Quantum mysteries for anybody, *Am. J. Phys.* **49** (10), 940–943 (1981).
- [9] N.D. Mermin, Quantum Mysteries for Anyone, *The Journal of Philosophy*, **78** (7) 397–408 (1981)
- [10] N.D. Mermin, Quantum mysteries revisited. *Am. J. Phys.* **58** (8), 731–734 (1990).
- [11] J.F. Clauser, M.A. Horne, A. Shimony, y R.A. Holt, Proposed experiment to test local hidden-variable theories, *Phys. Rev. Lett.*, **23** (15) 880–884 (1969)
- [12] J.F. Clauser y M.A. Horne, Experimental consequences of objective local theories, *Phys. Rev. D* **10** (2) (1974)
- [13] A. Aspect, P. Grangier y G. Roger, Experimental Tests of Realistic Local Theories via Bell's Theorem, *Phys. Rev. Lett.* **47** (7) (1981)
- [14] A. Aspect, J. Dalibard and G. Roger, Experimental Test of Bell's Inequalities Using Time-Varying Analyzers, *Phys. Rev. Lett.* **49** (25) 1804–1807
- [15] W. Tittel, J. Brendel, H. Zbinden y N. Gisin, Violation of Bell Inequalities by Photons More Than 10 km Apart, *Phys. Rev. Lett.* **81** (17) (1998)
- [16] A. Peres, What is a state vector?, *Am. J. Phys.* **52** 644–650 (1984).
- [17] F. Laloë, Do we really understand quantum mechanics? Strange correlations, paradoxes, and theorems, *Am. J. Phys.* **69** 655–701 (2001).
- [18] N.D. Mermin, What is quantum mechanics trying to tell us?, *Am. J. Phys.* **66** (9) 753 (1998)
- [19] No te preocupes si no entiendes, por ejemplo, la palabra “autoadjunto”: por regla general las palabras no tienen demasiada importancia. Sigue adelante quitando lo que no se entiende.
- [20] Existen excepciones a esta regla, pero son algo extrañas: si conozco el estado ψ del sistema, entonces puedo arreglar una medición, escogida en función de ψ , que no afecte el estado del sistema y que dé cierto resultado con probabilidad una. Otro caso especial es el de las llamadas “mediciones débiles” de las que se supone que perturban el sistema poco. Sin embargo, los experimentos que se realizan usándolas terminan arrojando resultados bien definidos en base a mediciones comunes y corrientes, y las predicciones de la teoría de mediciones débiles no difieren de la mecánica cuántica tradicional, de manera que no puede decirse que éstas contradicen la visión tradicional de la medición “fuerte”.
- [21] El trastorno, sin embargo, sí resultaría mayúsculo en el caso, que personalmente juzgo inverosímil, que se lleguen realmente a observar los elementos clásicos de realidad subyacentes a la mecánica cuántica. En este caso verdaderamente estaríamos capaces de actuar sobre sistemas distantes sin que el efecto disminuya con la distancia y de manera instantánea; muy poco del edificio de la física actual podría sobrevivir a esto.

Planetary migration in gaseous protoplanetary disks

Frédéric S. Masset, Instituto de Ciencias Físicas, UNAM

1. Introduction

The importance of the tidal interaction between a protoplanetary disk and a forming planet was first recognized long before the discovery of the first extrasolar planet in 1995. [16] discussed the case of Jupiter in a conservative protoplanetary nebula, and found that its semi-major axis should evolve as a result of the gravitational interaction between the planet and the nebula (although they could not determine whether it should increase or decrease).

When the first extrasolar planet was discovered orbiting 51 Peg with a period of 4,23 days () at a distance of only 0,052 AU from the central star, theories of orbital migration received renewed attention. None of the reasonable planetary formation scenarios was able to account for the formation of a planetary core that close to the star. It therefore appeared likely that this planet had formed farther out in the protoplanetary disk and then migrated towards the star, along the lines of predictions made by the theoretical work of the eighties ([31]).

Had any doubted that significant planetary migration is common in forming planetary systems, additional clues were provided by the discovery of planetary systems exhibiting low-order mean motion resonance. Under the effect of differential migration (i.e., the outer planet migrates inwards faster than the inner one), two planets can converge and be captured in a low-order mean motion resonance.

The present contribution is organized as follows: in section 2 we define the notation, in section 3, we present the torque expressions at the Lindblad and corotation resonances. In section 4, we present the different migration modes that have been envisaged so far. Finally, in section 5 we present a list of recent results that numerical simulations have recently brought to our knowledge of planet–disk interactions.

2. Notation

We consider a Keplerian gaseous disk with vertical scaleheight $H(r)$, surface density $\Sigma(r)$, kinematic viscosity $\nu(r)$, and orbital frequency $\Omega(r)$, where r is the distance to the central object. We consider a planet with a prograde orbit coplanar to the disk, of mass M_p and orbital frequency Ω_p . Whenever we consider a single azimuthal Fourier component of a given quantity, we denote by m its azimuthal wavenumber.

3. Disk torque at an isolated resonance

The problem of determining the torque between any perturbing potential and the disk, in the linear regime, amounts to determining the torque exerted on the disk by the Fourier components of

the potential. [15] have shown that angular momentum exchange between the perturbing potential and the disk occurs only at the Lindblad and corotation resonances. Lindblad resonances correspond to locations in the disk where the perturbing potential's frequency in the matter frame ($\tilde{\omega}(r) = m[\Omega_p - \Omega(r)]$) matches $\pm\kappa(r)$ (the epicyclic frequency). The corotation resonance occurs where the perturbing potential's frequency is zero in the matter frame, that is to say at a radius where the disk material rotates along with the perturbing potential.

3.1. Torque at a Lindblad resonance

3.1.1. Torque expression

The torque expression at a Lindblad resonance by a single Fourier component of the potential with m -fold symmetry is ([15], [44], [1, Artymowicz 1993])

$$\Gamma_m = -\frac{m\pi^2\Sigma}{rdD/dr} \left(r \frac{d\Phi_m}{dr} + \frac{2\Omega}{\Omega - \Omega_p} \Phi_m \right)^2, \quad (1)$$

where Γ_m is the torque exerted on the disk material by the perturbing potential, $D = \kappa(r)^2 - m^2[\Omega(r) - \Omega_p]^2$ represents a distance to the resonance and $\Phi_m(r)$ is the amplitude of the potential component. In Eq. (1), the term in brackets and rdD/dr are both to be evaluated at the resonance location. In a Keplerian disk, rdD/dr is positive at the ILR (Inner Lindblad Resonance, where $\tilde{\omega} = -\kappa$) and negative at the OLR (Outer Lindblad Resonance, where $\tilde{\omega} = +\kappa$). The perturbing potential therefore exerts a negative torque on the disk at the ILR, and a positive torque at the OLR. Newton's third law thus implies that the disk exerts a positive (negative) torque on the perturber at the ILR (OLR).

3.1.2. Lindblad resonance location

In order to evaluate the torques given by Eq. (1), one has to know the location of the Lindblad resonances. As stated previously, a Lindblad resonance is found where $\tilde{\omega} = \pm\kappa$ (the upper sign stands for the OLR, while the lower sign stands for the ILR). Using the fact that $\kappa = \Omega$ in a Keplerian disk, we obtain: $\Omega(r_{\text{LR}}) = \frac{m}{m\pm 1}\Omega_p$. Note that owing to pressure effects, the waves launched by the potential components are slightly offset from the resonance locations. In particular, as $m \rightarrow \infty$ the turning point locations tend to pile up at a radius given by: $r = r_c \pm \frac{\Omega}{2A}H$. These points of accumulation correspond to the radius at which the flow becomes supersonic in the corotating frame ([17]). In the case of a Keplerian disk, these points are located $\pm(2/3)H$ away from the corotation radius.

3.2. Torque at a corotation resonance

The angular momentum exchange at a corotation resonance and a Lindblad resonance are due to different physical processes. In the latter case the perturbing potential tends to excite epicyclic motion, and the angular momentum deposited is evacuated through pressure-supported waves. On the other hand, these waves are evanescent in the corotation region and therefore unable to remove the angular momentum brought there by the perturber ([15, Goldreich & Tremaine 1979]). In the linear regime, the corotation torque exerted by a perturbing potential with m -fold symmetry on the disk is

$$\Gamma_c = \frac{\pi^2 m}{2} \left[\frac{\Phi_m^2}{d\Omega/dr} \frac{d}{dr} \left(\frac{\Sigma}{B} \right) \right]_{r_c}, \quad (2)$$

where the term in brackets is to be evaluated at the corotation radius. The corotation torque is thus proportional to the gradient of Σ/B , evaluated at the corotation radius, and where B is equal to half the flow vorticity. The corotation torque is therefore proportional to the gradient of the vortensity (ratio of the vorticity to the surface density). The corotation torque is therefore zero in a disk with $\Sigma \propto r^{-3/2}$, such as the minimum mass solar nebula (MMSN).

The physical picture of the flow at a corotation resonance with azimuthal wavenumber m is characterized by a set of m eye-shaped libration islands in which fluid elements move along closed streamlines. The corotation torque is prone to saturation, which can be described as follows: when the disk viscosity is close to zero, the vortensity is conserved along a fluid element's path. The libration of fluid elements redistributes the vortensity within the libration islands. Once the vortensity has been sufficiently stirred up, even an infinitesimally small amount of viscosity suffices to render the vortensity uniform over the whole libration island. The corotation torque then goes to zero (i.e., saturates), because it scales with the vortensity gradient.

In order to avoid saturation, the viscosity must be high enough to prevent the vortensity from becoming uniform over the libration islands. This is possible if the viscous timescale across these islands is smaller than the libration timescale, as shown by [50]. In this case, viscous diffusion across the libration islands permanently imposes the large-scale vortensity gradient over the libration islands. Finally, it should be noted that saturation properties cannot be captured by a linear analysis, since saturation requires a finite libration time, and thus a finite resonance width.

4. Planetary migration

4.1. Type I migration

We consider the case of a low-mass planet, so that the overall disk response can be treated as a linear superposition of its responses to individual Fourier components of the potential. Each component torques the disk at its Lindblad and corotation resonances. We denote by Γ_{ILR}^m the torque of the m^{th} potential component at its ILR, and adopt similar notation for the torques at the Outer Lindblad Resonance (Γ_{OLR}^m) and corotation resonance (Γ_{CR}^m). The total tidal torque exerted by the disk on the planet, which is equal and opposite to the torque exerted by the planet on the disk, can therefore be written as

$$\Gamma = \sum_{m>0} \Gamma_{ILR}^m + \sum_{m>0} \Gamma_{OLR}^m + \sum_{m>0} \Gamma_{CR}^m. \quad (3)$$

The first series in this sum is the total Inner Lindblad torque, and the second is the total Outer Lindblad torque. The absolute value of either term is also called the one-sided Lindblad torque. The last term is called the coorbital corotation torque. The sum of the two Lindblad torques is referred to as the differential Lindblad torque.

4.1.1. Differential Lindblad torque

The acoustic shift of the effective Lindblad resonances mentioned at section 3 has an important consequence: there is a sharp cut-off in the high- m torque components (for $m \gg r/H$) as shown by [1, Artymowicz (1993)], since the high- m potential components become localized in increasingly narrow annuli around the perturber orbit. The value of a potential component at the accumulation point (where the torque is exerted) therefore tends to zero as m tends to infinity.

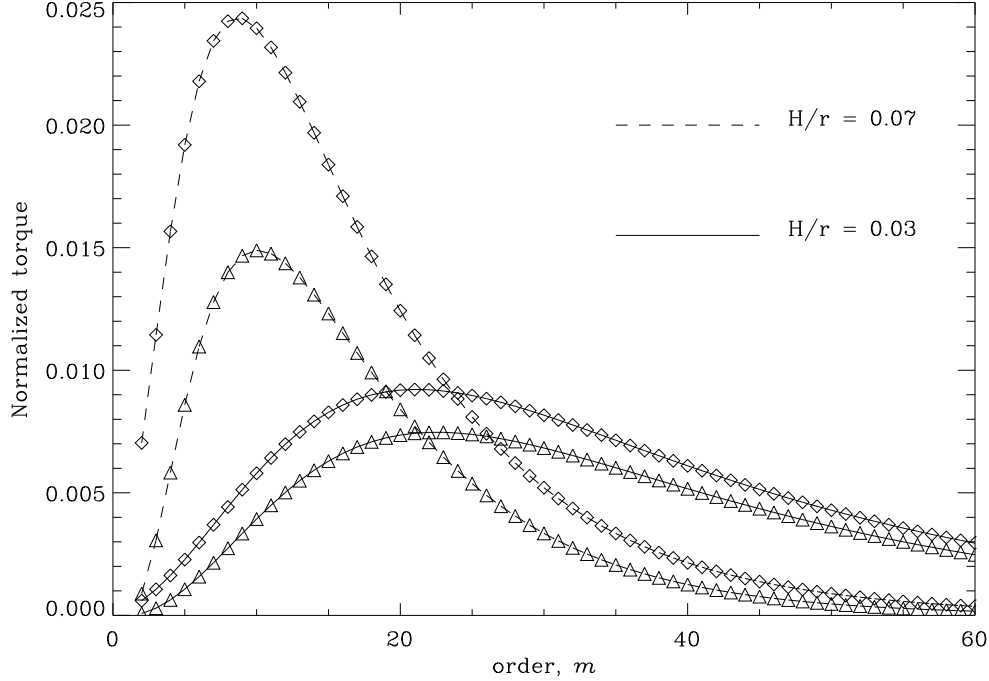


Figure 1: The absolute value of individual inner (triangle) and outer (diamond) torques as a function of m . The torques are normalized to the value $\Gamma_0 = \pi q^2 \Sigma a^4 \Omega_p^2 h^{-3}$, where q is the planet mass to star mass ratio. The one-sided Lindblad torques scale as $(H/r)^{-3}$, hence the total areas under each curve are of the same order of magnitude. The differential Lindblad torque scales as $(H/r)^{-2}$.

Fig. 1 illustrates the behavior the one-sided Lindblad torques. In particular, one can see that the cutoff occurs at larger m in a thinner disk. Also, for both disk aspect ratios there is a very apparent mismatch between the inner and the outer torques; the former is systematically smaller than the later. If we consider the torque of the disk acting on the planet, then the outer torque is negative and the inner torque is positive; the total torque on the planet is therefore negative. As a consequence, migration is directed inwards and the orbit decays towards the central object ([65]).

4.1.2. Pressure buffer

One remarkable feature of the differential Lindblad torque is its weak dependence on the slope of the surface density function. This is not what one would naively expect, since as one increases the slope the surface density increases at the Inner Lindblad Resonances and decreases at the Outer Lindblad Resonances. As one increases the surface density gradient, however, one simultaneously increases the radial pressure gradient. This makes the disk more and more sub-Keplerian. As a consequence, the Outer Lindblad Resonances approach the planet's orbit while the Inner Lindblad Resonances recede from it. This process plays against the more obvious effect of the surface density. This effect is known as the pressure buffer ([66], [62]), and frustrates any reasonable attempt to revert the differential Lindblad torque by tuning the power law indexes of the surface density and temperature profiles. This makes inward type I migration inevitable, at least in disks where the surface density and temperature are power laws of the radius.

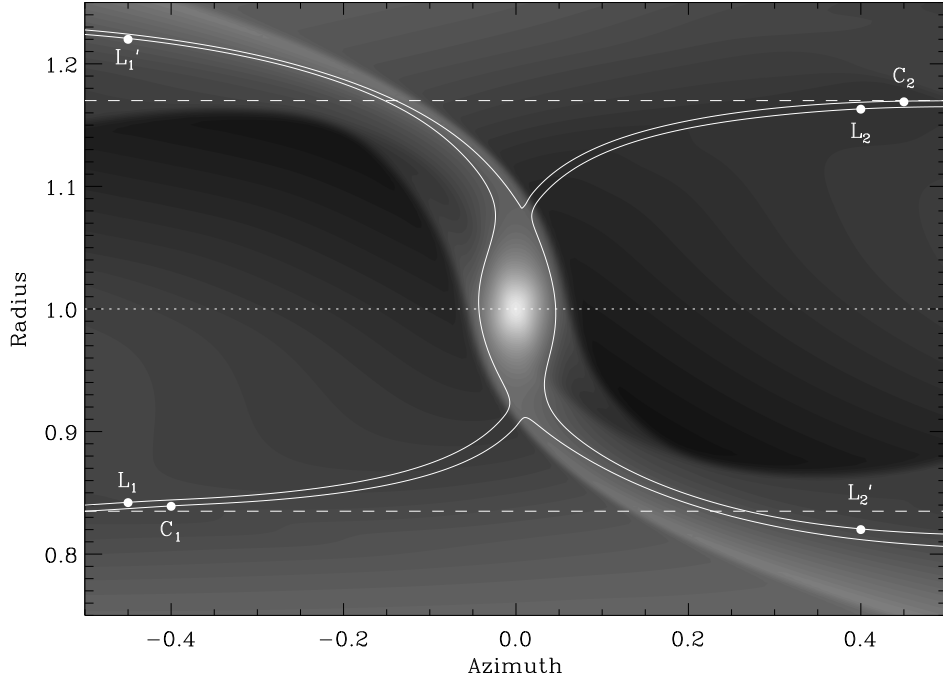


Figura 2: Asymmetry of the horseshoe region. The circulating fluid elements C_2 (moving towards the left) and C_1 (moving towards the right) recede from the orbit after crossing the shock excited by the planet. Similarly, the librating fluid elements recede from the orbit after executing their horseshoe U-turns. This particular example shows streamlines of the flow in the corotating frame of a $2 M_J$ planet in a disk with $h = 0,05$. The planet is on a fixed circular orbit, and this snapshot was taken after 22,5 orbits.

4.1.3. Type I migration timescale

The most up-to-date estimate of the total (*i.e.* Lindblad plus corotation) tidal torque in the linear regime between a three dimensional disk and a low mass planet is the estimate by [62]. It yields a migration timescale of 8×10^5 yrs for an Earth-mass object embedded at 5 AU in the MMSN. This is much shorter than the disk lifetime.

4.2. Type II migration

4.2.1. Shock appearance and horseshoe asymmetry

The wake excited by a planet eventually turns into a shock. The location at which profile steepening produces a shock depends on the planet mass; the larger the mass, the closer the shock will be to the orbit. For planets above some critical mass, the wake becomes a shock within the excitation region. Under these circumstances, the fluid elements circulating just outside the co-orbital region receive a kick of angular momentum every time they cross the wake. This is represented in Fig. 2. As a consequence horseshoe U-turns are not symmetric. A fluid element initially located inside the libration region thus progressively recedes from the orbit as it performs a sequence of horseshoe U-turns, until it ends up in the inner disk or the outer disk ([35]). The co-orbital region is thereby

emptied, and an annular gap eventually appears around the orbit. The timescale for emptying the co-orbital region can readily be estimated from Fig. 2. After each horseshoe U-turn, the distance of a fluid element from the orbit increases by an amount between 10 and 20%. The characteristic emptying time of the horseshoe region is therefore $O(10)$ times half the libration time, which is given by $\tau_{\text{lib}}/2 = 2\pi a/(3/2)\Omega_p x_s = (2/3)T_o(a/x_s)$. Here we can estimate from the figure that $x_s \approx 0,16$, so $\tau_{\text{lib}}/2 \approx 4 T_0$. In this particular example, the co-orbital region is therefore emptied after about $O(40)$ orbits. This simple estimate also shows that the smaller the planet mass, the longer the gap clearance timescale. Indeed, as the planet mass decreases, the horseshoe region becomes more and more symmetric so that more libration times are needed to get rid of the co-orbital material, while the libration time itself increases. For a $1 M_J$ planet orbiting in a disk with $H/r = 0,05$, the clearance timescale of the gap is about 100 orbits.

4.2.2. Accretion

A planet engaged in type II migration has a mass much larger than the critical mass for runaway gas accretion ([60]). It therefore accretes gas from the nebula at the same time as it migrates. [25] has devised a scheme to simulate the accretion in the planet’s vicinity, and found that a giant planet can still accrete significant amounts of gas despite the presence of the gap. [49] have found that the final mass of a protogiant of initially one Jupiter mass could be of several Jupiter masses. The ability of a giant planet to accrete the surrounding gas depends on the equation of state of the gas and its ability to get rid of the gravitational energy released by the accretion. [9], [24] and [61] have performed high resolution hydrodynamics calculations taking into account radiative transfer in order to assess the dynamics of the inner Roche lobe and its impact on accretion. The former still find significant accretion while this last work shows that the accretion is stalled. In all cases the geometry of the flow in the Roche lobe is that of a bubble rather than a thin accretion disk.

4.2.3. Gap opening criteria

Classically, the gap opening conditions once consisted of two independent criteria ([32], [34], [7]) that needed to be simultaneously fulfilled. The first, referred to as the thermal criterion (since it imposes a limit on the disk thickness, and hence on the disk temperature) requires that the wake becomes a shock just as it is excited. The flow must therefore be strongly non-linear in the planet’s vicinity, and the parameter R_H/H must be larger than some critical value, where R_H is the planetary Hill radius. This critical value is ~ 1 , although its precise value can be slightly different. The second criterion is that the viscosity is sufficiently low, so that the surface density jump across the edges of the excavated region is a sizable fraction of the unperturbed surface density. This condition, which is known as the viscous criterion, is expressed as $q > \frac{40}{\mathcal{R}}$ where $\mathcal{R} = a^2\Omega_p/\nu$ is the Reynolds number. [8] have used another condition, namely that the circulating streamline just outside the separatrix should be closed, to derive the gap surface density profile semi-analytically. They require that the integral of the viscous, gravitational, and pressure torques cancels out over one synodic period of a given fluid element. They provide an ansatz expression for the pressure torque that is approximately valid for a reasonable range of planetary masses and disk thicknesses, and derive the following unique criterion for gap opening: $\frac{3}{4}\frac{H}{R_H} + \frac{50}{q\mathcal{R}} < 1$. Broadly speaking, this criterion is approximately equivalent to the previous two except in the case where both are only marginally fulfilled.

4.2.4. Migration of planets that open a gap

A “clean” gap (i.e., a gap with little residual surface density) splits the disk material into an outer disk and an inner disk. Therefore, the planet must drift inwards at the same rate that the outer disk spreads inwards. In other words, the migration rate of a giant planet that has opened a gap in the disk is the same as the viscous drift rate of the disk ([33]). This type of migration is referred to as type II migration ([49]). It is usually said that in this regime, the planet’s orbit is locked to the disk’s viscous evolution. The migration drift rate of the planet is therefore

$$\frac{da}{dt} \sim -\frac{v}{a}. \quad (4)$$

For a $M_p = 1 M_J$ planet that undergoes type II migration in a disk with $H/r = 0,04$ and $\alpha = 6 \cdot 10^{-3}$, the migration time starting from $a = 5$ AU is about $1,6 \cdot 10^4$ orbits. This corresponds to $\sim 1,6 \cdot 10^5$ years, if the central object has one solar mass.

Using two-dimensional numerical simulations, [49] have shown that the migration of giant planets (with masses greater than or equal to one Jupiter mass) in a viscous disk obeys the scenario outlined above, at least broadly speaking. In particular, they found that the timescale of variation in the planet’s semi-major axis is similar to the viscous timescale of the disk. These results have been obtained by assuming that the effective viscosity of the disk is adequately modeled by the Navier–Stokes equation. In this approach the kinematic viscosity is chosen to account for the accretion rates inferred from observations of T Tauri objects. [56] have performed much more numerically demanding calculations; instead of resorting to the purely hydrodynamical scheme including an *ad hoc* kinematic viscosity, their model describes the self-sustained magnetohydrodynamic (MHD) turbulence arising from the magnetorotational instability (MRI).

They find that a giant protoplanet still opens a gap in the disk, in much the same manner as in a disk modeled by the Navier–Stokes equations. Surprisingly, the gap in a turbulent disk tends to be larger and deeper than in a laminar disk ([55]). The mass accretion rate tends to be larger in the MHD turbulent case, most likely because of magnetic breaking of the circumplanetary disk ([56]).

4.2.5. Type II migration of several planets

The migration properties of two or more giant planets is a topic that has received a lot of attention, primarily because we detect extrasolar giant planets that are in mean motion resonance, which is a natural outcome of convergent type II migration ([36], [28], [27]), and secondly because under some circumstances (if the outer planet is sufficiently lightweight and barely opens a gap), the migration of the whole system may be reversed and be directed outwards ([36]). This scenario, applied to the case of our own Solar System, has received the name of Grand Tack, and has been invoked to explain a number of properties of the asteroid belts and planets [64]. Finally, as was first noted by [26] the distance between the giant planets brought to close orbits by convergent migration can be sufficiently short to render the system unstable after gas clearance, which may account for the eccentric orbits of some extrasolar planets.

4.3. Type III migration

Type III migration refers to a mode of migration for which the major driver is material flowing through the coorbital region. In the previous sections, the torque acting on a migrating planet was

considered independent of its migration rate. However, the corotation torque implies material that crosses the planet orbit on the U-turn of the horseshoe streamlines. In a non-migrating case, only the material trapped in the horseshoe region participates in these U-turns, but in the case of an inward (or outward) migrating planet, material of the inner disk (outer disk) has to flow across the co-orbital region and executes one horseshoe U-turn to do so. By doing this, it exerts a corotation torque on the planet that scales with the drift rate. We call x_s the half radial width of the horseshoe region. The amount of specific angular momentum that a fluid element near the separatrix takes from the planet when it crosses the planet orbit and goes from the orbital radius $a - x_s$ to the orbital radius $a + x_s$ is $\Omega_p a x_s$. The corresponding torque exerted on the planet in steady migration is therefore, to lowest order in x_s/a :

$$\Gamma_2 = (2\pi a \Sigma_s \dot{a}) \cdot (\Omega_p a x_s), \quad (5)$$

where we keep the same notation as in [41], and where Σ_s is the surface density at the upstream separatrix. As the system of interest, we take the system composed of the planet and all fluid elements trapped in libration in its co-orbital region, namely the whole horseshoe region (with mass M_{HS}) and the Roche lobe content (with mass M_R), because all of these parts perform a simultaneous migration. The drift rate of this system is then given by :

$$(M_p + M_{HS} + M_R) \cdot (a \dot{a} \Omega_p / 2) = (4\pi a x_s \Sigma_s) \cdot (a \dot{a} \Omega_p / 2) + \Gamma_{LR} \quad (6)$$

which can be rewritten as:

$$m_p \cdot (a \dot{a} \Omega_p / 2) = (4\pi a \Sigma_s x_s - M_{HS}) \cdot (a \dot{a} \Omega_p / 2) + \Gamma_{LR} \quad (7)$$

where $m_p = M_p + M_R$ is all the mass content within the Roche lobe, which for now on for convenience we refer to as the planet mass. The first term of the first bracket of the r.h.s. corresponds to the horseshoe region surface multiplied by the upstream separatrix surface density, hence it is the mass that the horseshoe region would have if it had a uniform surface density equal to the upstream surface density. The second term is the actual mass of the horseshoe region. The difference between these two terms is called in MP03 the coorbital mass deficit and denoted δm . Eq (7) yields a drift rate :

$$\dot{a} = \frac{\Gamma_{LR}}{2Ba(m_p - \delta m)} \quad (8)$$

This drift rate is faster than the standard estimate in which one neglects δm . This comes from the fact that the coorbital dynamics alleviates the task of the differential Lindblad torque by advecting fluid elements from the upstream to the downstream separatrix. The angular momentum they extract from the planet by doing so favors its migration. As δm tends to m_p , most of the angular momentum lost by the planet and its coorbital region is gained by the orbit crossing circulating material, making migration increasingly cost effective. When $\delta m \geq m_p$, the above analysis, assuming a steady migration (\dot{a} constant), is no longer valid. Migration undergoes a runaway, and has a strongly time varying migration rate, that increases exponentially over the first libration times. Runaway (also said type III) migration is therefore a mode of migration of planets that deplete their coorbital region and embedded in sufficiently massive disks, so that the above criterion be satisfied. An analysis similar to the above calculation may be performed, in which the corotation torque depends on the migration rate, except that one now has to introduce a delay τ between the mass inflow at the upstream separatrix and the corotation torque. Fluid elements passing through the upstream separatrix need indeed on average

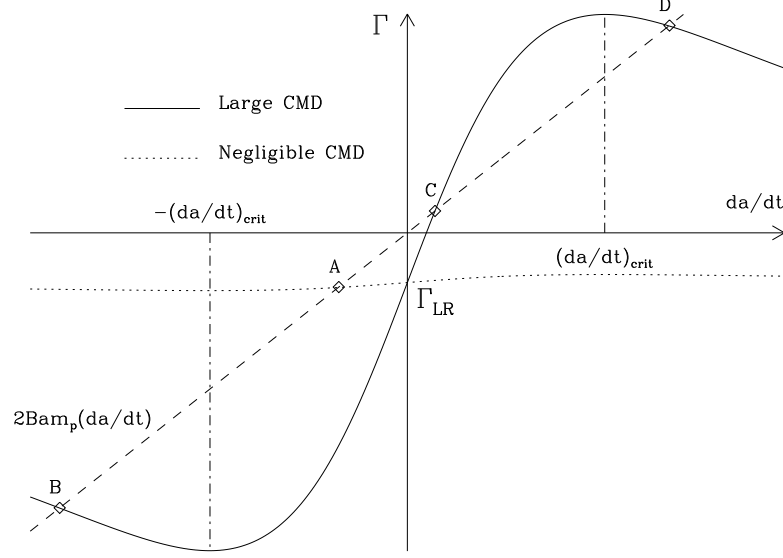


Figure 3: The solid curve shows the total torque on the planet in a massive disk (hence with a large coorbital mass deficit) as a function of the drift rate. For $|\dot{a}| \ll \dot{a}_{\text{crit}}$ the torque exhibits a linear dependence in \dot{a} . The dotted line shows the torque in a low mass disk (i.e. with a negligible coorbital mass deficit), in which case the torque is almost independent of the migration rate and is always close to the differential Lindblad torque Γ_{LR} . The dashed line represents the planet angular momentum gain rate as a function of \dot{a} , assuming a circular orbit. For a given situation, the migration rate achieved by a steadily migrating planet is given by the intersection of the dashed line with the torque curve. In the low mass disk case, the intersection point, A, is unique, and stable. It yields a negative drift rate controlled by the differential Lindblad torque. In the high mass disk case (type III case), there are 3 points of intersection (B, C and D). The central point (C) is unstable, while the extreme ones (B and D) are stable and correspond to the maximum drift attained by the planet, either inwards (point B) or outwards (point D).

a fraction of a libration timescale to reach the planet and execute a horseshoe U-turn. This delay represents the latency of the feedback loop.

$$\Gamma_{CR}(t) = 2Ba\delta m\dot{a}(t - \tau) \quad (9)$$

A Taylor expansion in time of $\dot{a}(t - \tau)$ yields a first order differential equation for \dot{a} (see MP03 for details). The linear dependence of the corotation on the drift rate remains valid as long as the semi-major axis variation over a horseshoe libration time is smaller than the horseshoe zone width, i.e.:

$$|\dot{a}| < \dot{a}_{\text{crit}} = \frac{Ax_s^2}{2\pi a} \quad (10)$$

The corotation torque then reaches a maximum and slowly decays for larger values of \dot{a} (see fig. 3). The terminal drift rate of a type III steadily migrating planet can be estimated by a standard bifurcation analysis as illustrated in fig. 3. The transition from one case to the other (one intersection point to three intersection points) occurs when the line showing the rate of change of the angular momentum

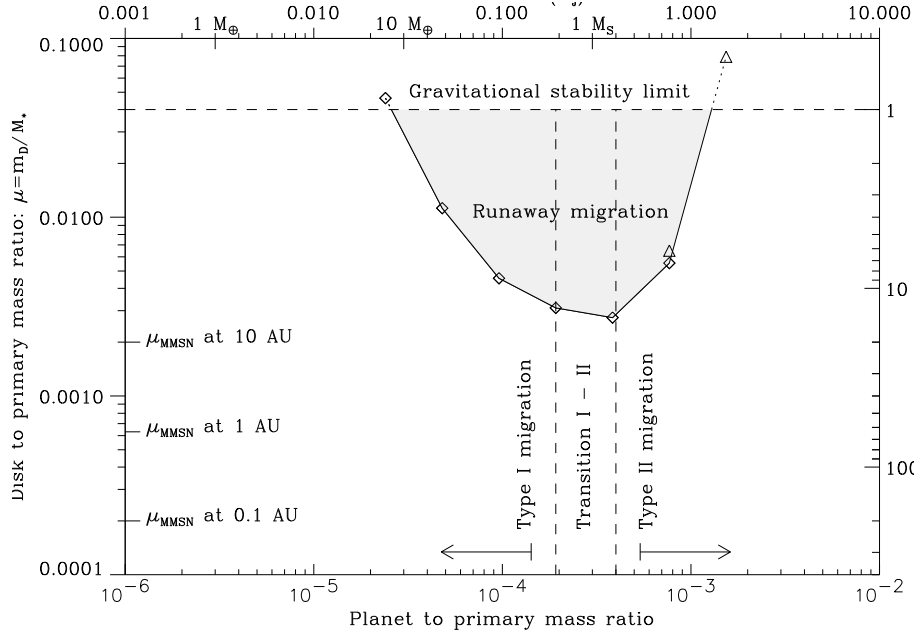


Figure 4: Runaway (or type III) limit domain for a $H/r = 0,04$ and $\nu = 10^{-5}$ disk, with a surface density profile $\Sigma \propto r^{-3/2}$. The variable $m_D = \pi\Sigma r^2$ features on the y axis. It is meant to represent the local disk mass, and it therefore depends on the radius. Type III is most likely for Saturn mass planets. These would undergo type III migration in disks no more massive than a few times the MMSN.

(which has slope $am_p\Omega_p/2$) and the torque curve near the origin are parallel. Since this latter has a slope $a\delta m\Omega_p/2$ near the origin, the transition occurs near $m_p = \delta m$. The disk critical mass above which a planet of given mass undergoes a runaway depends on the disk parameters (aspect ratio and effective viscosity). The limit has been worked out by MP03 for different disk aspect ratios and a kinematic viscosity $\nu = 10^{-5}$. We reproduce in figure 4 the type III migration domain for a disk with $H/r = 0,04$. A number of comments can be made from figure 4:

- The MMSN was barely massive enough to yield type III migration of Saturn. This suggests that in many protoplanetary disks, inferred to be several times more massive than the MMSN, type III migration is very likely for Saturn mass protoplanets.
- Type III is impossible for massive planets ($m_p > 1 M_J$) as the horseshoe separatrices sample the gap edges in regions significantly depleted, yielding a small coorbital mass deficit.
- The sharp limit on the high mass side of the runaway domain might be related to the fact that most of the extrasolar planets known as “hot Jupiters”, with a semi-major axis $a < 0,2$ AU, happen to have sub-Jovian masses. A forming protoplanet, as it passes through the runaway domain, would migrate very fast towards the central object in a series of type III episodes, and at the same time it would accrete gas from the nebula. If the protoplanet happens to get out of the runaway domain before it reaches the central regions of the disk, it enters the slow, type II migration regime, having at least about a Jupiter mass. Otherwise, it may reach the central regions through type III migration (if the surface density profile is steep enough), still as a sub-Jovian object.

Type III migration, for the same disk profile and planet mass, can be directed either outwards or inwards, depending on the initial conditions. This is related to the fact that the differential equation obeyed by the semi-major axis, in the type III regime, is a second order differential equation in which one can specify independently a and \dot{a} . The type III migration of a planet therefore depends on its migration history, the “memory” of this history being stored in the way the horseshoe region is populated, i.e. in the repartition of the coorbital mass deficit. Note that owing to the strong variation of the drift rate in a type III episode, the horseshoe streamlines are not exactly closed, so that the coorbital mass deficit can be lost and type III migration can stall. This has been observed in numerical simulations, which show that the semi-major axis is varied by a factor of a few at most during a single type III episode. Since the early work of [41], the effect of migration on the coorbital dynamics has been extended to lower mass planets by [51].

4.4. Stochastic migration

Oftentimes, the protoplanetary disk is considered laminar while the accretion of material onto the central object is ensured by an *ad hoc* kinematic viscosity chosen to account for the mass accretion rate measured for T Tauri stars. The molecular viscosity of protoplanetary disks appears to be insufficient by many orders of magnitude, however, to reproduce the accretion rates typically measured. The source of the high effective viscosity in these disks is thought to be turbulence. The MRI (see section 4.2) has been identified as a powerful source of MHD turbulence in magnetized disks ([2, Balbus & Hawley 1991], [19, Hawley & Balbus 1991], [20, Hawley & Balbus 1992]), and this section will exclusively focus on the impact of this kind of turbulence on planetary migration.

MRI can develop only in regions of the disk where the matter and magnetic field are coupled, which requires a sufficiently high (albeit weak) ionization rate. In the planet-forming region (1 – 10 AU), it is thought that only the upper layers of the disk are ionized by X-rays from the central star or cosmic rays ([14], [12]). The bulk of the disk, however, should be ionized outside this region. This has led [14] to the concept of layered accretion: the upper layers of the region between 1 and 10 AU participate in accretion onto the central star, whereas its magnetically inactive equatorial parts, usually called the *dead zone*, do not participate in the inwards flow of disk material.

There already exist a large number of works describing numerical simulations that self-consistently describe an MHD turbulent disk with embedded planets ([48], [57], [47], [46], [67]). They exclusively consider a fully magnetized disk (hence with no dead zone), however, without any vertical stratification for reasons of computational cost. More recently, [3] have examined more specifically the behavior of the horseshoe drag in a disk invaded by MHD turbulence. They found that this component of the torque is still active, and that fluid elements still execute horseshoe U-turns in the vicinity of the planet. [18] have contemplated the horseshoe dynamics in a 2D disk threaded by a toroidal magnetic field and found the latter to have a considerable impact on the net torque, when the disk’s viscosity and resistivity are taken into account.

Not surprisingly, the torque felt by a planet in a turbulent disk displays large temporal fluctuations. One can assign an order of magnitude to their amplitude by considering an overdense region of size H , located at a distance H from the planet such that the perturbed density in this region is of the same order as the unperturbed density. This yields an order of magnitude for the torque fluctuations of $G\Sigma a$ ([48, Nelson & Papaloizou 2004], [46, Nelson 2005]).

[48, Nelson and Papaloizou (2004)] and [46, Nelson (2005)] have investigated the migration of low and intermediate mass planets embedded in turbulent disks. [30, Laughlin et al. (2004)] have also investigated this problem, but rather than tackling it through self-consistent numerical simulations

they performed a two-dimensional calculation which mimicked the effects of turbulence using a time-varying, non-axisymmetric potential acting on the gas disk, rather than directly on the planet. The migration of low-mass planets embedded in turbulent disks is significantly different from the type I migration expected for laminar disks. The large torque variations due to turbulence induce the planet’s semi-major axis to evolve on a random walk rather than systematically decay.

One question that is still open is whether the total torque felt by a planet in a turbulent disk can be decomposed into a laminar torque and the effect of fluctuations arising from the turbulence. We call the latter component the stochastic torque. One might expect that the time average of the stochastic torques is negligible compared to the total mean torque (which might be the same as the laminar torque, but this is still unknown), provided that this average is performed over a time interval that is much longer than the turbulence recurrence time. Under this assumption, the behavior of the planet should exhibit a systematic trend reminiscent of type I migration. [46, Nelson (2005)] has investigated the statistical properties of these torque fluctuations, finding significant power at low frequencies, corresponding to timescales comparable to the simulation time. As a consequence, in many of his calculations no systematic trend is observed; stochastic migration dominates type I migration over the entire run time of his calculations, or about 150 orbits. The reason for such significant power at very low frequencies is still unknown.

The amplitude of the specific stochastic torque is independent of the planet mass, whereas the specific wake torque scales with the planet mass. [46, Nelson (2005)] found that for planets up to $\sim 10 M_{\oplus}$ the stochastic migration overcomes the systematic trend (over a simulation run time of 150 orbits), whereas systematic effects are dominant for larger masses. We mention however the recent work by [11, Fromang & Nelson (2006)], who argue that density fluctuations are smaller in a stratified, turbulent disk than in the unstratified models currently used to assess stochastic torques. This argument suggests that systematic effects could be dominant at masses even lower than $10 M_{\oplus}$.

As pointed out by [23, Johnson et al. (2006)], if the turbulence has a finite correlation time then the stochastic (or diffusive) migration of low-mass planets can be reduced to an advection-diffusion equation. They show that diffusion always reduces the mean migration time of the planets, although a fraction of them still “survive” an extended period of migration.

5. Additional results on planetary migration

The use of ever increasing computational resources, and the inclusion of new physical ingredients have yielded over the years a number of new results which have somehow modified the classical textbook picture drawn in the first sections of this chapter. We draw hereafter a non comprehensive list of such results.

5.1. Planetary migration and magnetic field

Notwithstanding the issue of MRI and its non-linear outcome as MHD turbulence, the role of a toroidal or poloidal magnetic field on type I migration in a laminar disk as been contemplated by several authors. [63] considers a disk threaded by a toroidal magnetic field, and shows that when the magnetic field as a function of radius decreases sufficiently fast, the total torque felt by the planet is positive, hence the planet migrates outwards. [13] have performed two-dimensional numerical simulations which essentially confirmed the analytic predictions of [63]. More recently, Muto et al. (these

proceedings) worked out the analytic torque expression both for a disk threaded by a poloidal magnetic field and a disk with a toroidal magnetic field, which enables them to make a variety of predictions about type I migration in magnetized disks.

5.2. Inclusion of the disk's self gravity

The protoplanetary disk's self-gravity, usually neglected on the grounds of its large Toomre parameter, has recently been contemplated by a number of authors, which were led to contradictory statements ([45], [58]). [5] have revisited previous works on this subject and found that the inclusion of self-gravity slightly speeds up migration with respect to analytical drift rate estimates. They also exhibited a strong bias that systematically affects numerical calculations in which a planet is released and freely migrates in a non self-gravitating disk.

5.3. Role of the corotation torque at a cavity edge

Eq. (2) shows that the corotation torque can be a large, positive quantity at a surface density jump such that the surface density is larger on the outside, and that it may possibly overcome the differential Lindblad torque. This happens at a cavity edge, even if the cavity is shallow. This has led [40] to the concept of planetary trap: type I migrating embryos are stopped whenever they reach a relatively abrupt drop of surface density, such as could be found at the inner edge of a dead zone, at the inner edge of a tidally truncated disk [59] or at the snow line [29].

5.4. Planetary migration and radiative transfer

Radiative transfer plays a very important role for planetary migration scenarios, for many different reasons. [43] exploit the differential Lindblad torque's extreme sensitivity to the location of the Lindblad resonances. They consider realistic models of T Tauri α -disks instead of the customary power law models, and show that type I migration can be significantly slowed at opacity transitions. [21] argue that taking into account the temperature perturbations due to shadowing and irradiation of the disk photosphere could significantly reduce the type I migration rate. [53] (hereafter PM06) consider a low-mass planet embedded in a disk with inefficient radiative cooling. A complex temperature structure develops in the vicinity of the planet which gives an underdense region behind the planet. As a consequence, the disk ultimately exerts a positive torque on the planet. This result clearly indicates that radiative transfer effects may prove crucial in resolving the problem of type I drift.

[4], [54], [38], [39] [52] [22] have undertaken a follow up study of the results of PM06 in order to provide ever increasing accurate formulae of the tidal torque experienced by a low mass planet in a gaseous protoplanetary disk.

6. Conclusion

Significant progress has been recently accomplished in the theories of planet-disk tidal interactions. Most of the new results have primarily been brought by large-scale calculations using modern supercomputer resources. In particular, the problem of type I migration has eventually undergone a considerable shift of perspective, by relaxing the customary barotropic assumption, thereby enabling a

new kind of perturbation (the so-called contact discontinuity familiar to the Riemann solvers community) to arise in the co-orbital region in the presence of an entropy gradient. In a different vein, recent work has highlighted the importance of the disc disturbance that does not arise from the embryo's gravity, but from the luminosity due to accretion [6, 10, 37, 42]. These effects can revert migration and excite eccentricity and inclination, with potentially fundamental consequences on scenarios of planetary formation.

The actual trend among numericists performing calculations of planet-disk interactions is to include more and more physics relevant to planetary migration in their schemes. These efforts render scenarios of planetary migration progressively more quantitative and predictive, and they should in the future eventually bridge the gap between the properties of the protoplanetary disk, and the structure of the planetary systems that may emerge from it.

Referencias

- [1] P. Artymowicz. On the Wave Excitation and a Generalized Torque Formula for Lindblad Resonances Excited by External Potential. *Astrophys. J.*, 419:155–+, December 1993.
- [2] S. A. Balbus and J. F. Hawley. A powerful local shear instability in weakly magnetized disks. I - Linear analysis. II - Nonlinear evolution. *Astrophys. J.*, 376:214–233, July 1991.
- [3] C. Baruteau, S. Fromang, R. P. Nelson, and F. Masset. Corotation torques experienced by planets embedded in weakly magnetized turbulent discs. *Astron. Astrophys.*, 533:A84+, September 2011.
- [4] C. Baruteau and F. Masset. On the Corotation Torque in a Radiatively Inefficient Disk. *Astrophys. J.*, 672:1054–1067, January 2008.
- [5] C. Baruteau and F. Masset. Type I Planetary Migration in a Self-Gravitating Disk. *Astrophys. J.*, 678:483–497, May 2008.
- [6] P. Benitez-Llambay, F. Masset, G. Koenigsberger, and J. Szulagyi. Planet heating prevents inward migration of planetary cores. *Nature*, 520:63–65, April 2015.
- [7] G. Bryden, X. Chen, D. N. C. Lin, R. P. Nelson, and J. C. B. Papaloizou. Tidally Induced Gap Formation in Protostellar Disks: Gap Clearing and Suppression of Protoplanetary Growth. *Astrophys. J.*, 514:344–367, March 1999.
- [8] A. Crida, A. Morbidelli, and F. Masset. On the width and shape of gaps in protoplanetary disks. *Icarus*, 181:587–604, April 2006.
- [9] G. D'Angelo, T. Henning, and W. Kley. Thermohydrodynamics of Circumstellar Disks with High-Mass Planets. *Astrophys. J.*, 599:548–576, December 2003.
- [10] H. Eklund and F. S. Masset. Evolution of eccentricity and inclination of hot protoplanets embedded in radiative discs. *Mon. Not. R. Astron. Soc.*, 469:206–217, July 2017.
- [11] S. Fromang and R. P. Nelson. Global MHD simulations of stratified and turbulent protoplanetary discs. I. Model properties. *Astron. Astrophys.*, 457:343–358, October 2006.

- [12] S. Fromang, C. Terquem, and S. A. Balbus. The ionization fraction in α models of protoplanetary discs. *Mon. Not. R. Astron. Soc.*, 329:18–28, January 2002.
- [13] S. Fromang, C. Terquem, and R. P. Nelson. Numerical simulations of type I planetary migration in non-turbulent magnetized discs. *Mon. Not. R. Astron. Soc.*, 363:943–953, November 2005.
- [14] C. F. Gammie. Layered Accretion in T Tauri Disks. *Astrophys. J.*, 457:355–+, January 1996.
- [15] P. Goldreich and S. Tremaine. The excitation of density waves at the Lindblad and corotation resonances by an external potential. *Astrophys. J.*, 233:857–871, November 1979.
- [16] P. Goldreich and S. Tremaine. Disk-satellite interactions. *Astrophys. J.*, 241:425–441, October 1980.
- [17] J. Goodman and R. R. Rafikov. Planetary Torques as the Viscosity of Protoplanetary Disks. *Astrophys. J.*, 552:793–802, May 2001.
- [18] J. Guilet, C. Baruteau, and J. C. B. Papaloizou. Type I planet migration in weakly magnetized laminar discs. *Mon. Not. R. Astron. Soc.*, 430:1764–1783, April 2013.
- [19] J. F. Hawley and S. A. Balbus. A Powerful Local Shear Instability in Weakly Magnetized Disks. II. Nonlinear Evolution. *Astrophys. J.*, 376:223–+, July 1991.
- [20] J. F. Hawley and S. A. Balbus. Three-Dimensional Simulations of the MHD Shearing Instability in Accretion Disks. *Bulletin of the American Astronomical Society*, 24:1234–+, September 1992.
- [21] H. Jang-Condell and D. D. Sasselov. Type I Migration in a Nonisothermal Protoplanetary Disk. *Astrophys. J.*, 619:1123–1131, February 2005.
- [22] M. A. Jiménez and F. S. Masset. Improved torque formula for low and intermediate mass planetary migration. *ArXiv e-prints*, July 2017.
- [23] E. T. Johnson, J. Goodman, and K. Menou. Diffusive Migration of Low-Mass Protoplanets in Turbulent Disks. *Astrophys. J.*, 647:1413–1425, August 2006.
- [24] H. Klahr and W. Kley. 3D-radiation hydro simulations of disk-planet interactions. I. Numerical algorithm and test cases. *Astron. Astrophys.*, 445:747–758, January 2006.
- [25] W. Kley. Mass flow and accretion through gaps in accretion discs. *Mon. Not. R. Astron. Soc.*, 303:696–710, March 1999.
- [26] W. Kley. On the migration of a system of protoplanets. *Mon. Not. R. Astron. Soc.*, 313:L47–L51, April 2000.
- [27] W. Kley, M. H. Lee, N. Murray, and S. J. Peale. Modeling the resonant planetary system GJ 876. *Astron. Astrophys.*, 437:727–742, July 2005.
- [28] W. Kley, J. Peitz, and G. Bryden. Evolution of planetary systems in resonance. *Astron. Astrophys.*, 414:735–747, February 2004.

- [29] K. A. Kretke and D. N. C. Lin. Grain Retention and Formation of Planetesimals near the Snow Line in MRI-driven Turbulent Protoplanetary Disks. *Astrophys. J. Lett.*, 664:L55–L58, July 2007.
- [30] G. Laughlin, A. Steinacker, and F. C. Adams. Type I Planetary Migration with MHD Turbulence. *Astrophys. J.*, 608:489–496, June 2004.
- [31] D. N. C. Lin, P. Bodenheimer, and D. C. Richardson. Orbital migration of the planetary companion of 51 Pegasi to its present location. *Nature*, 380:606–607, April 1996.
- [32] D. N. C. Lin and J. Papaloizou. Tidal torques on accretion discs in binary systems with extreme mass ratios. *Mon. Not. R. Astron. Soc.*, 186:799–812, March 1979.
- [33] D. N. C. Lin and J. Papaloizou. On the tidal interaction between protoplanets and the protoplanetary disk. III - Orbital migration of protoplanets. *Astrophys. J.*, 309:846–857, October 1986.
- [34] D. N. C. Lin and J. C. B. Papaloizou. On the tidal interaction between protostellar disks and companions. In E. H. Levy and J. I. Lunine, editors, *Protostars and Planets III*, pages 749–835, 1993.
- [35] S. H. Lubow, M. Seibert, and P. Artymowicz. Disk Accretion onto High-Mass Planets. *Astrophys. J.*, 526:1001–1012, December 1999.
- [36] F. Masset and M. Snellgrove. Reversing type II migration: resonance trapping of a lighter giant protoplanet. *Mon. Not. R. Astron. Soc.*, 320:L55+, February 2001.
- [37] F. S. Masset. Coorbital thermal torques on low-mass protoplanets. *Mon. Not. R. Astron. Soc.*, 472:4204–4219, August 2017.
- [38] F. S. Masset and J. Casoli. On the Horseshoe Drag of a Low-Mass Planet. II. Migration in Adiabatic Disks. *Astrophys. J.*, 703:857–876, September 2009.
- [39] F. S. Masset and J. Casoli. Saturated Torque Formula for Planetary Migration in Viscous Disks with Thermal Diffusion: Recipe for Protoplanet Population Synthesis. *Astrophys. J.*, 723:1393–1417, November 2010.
- [40] F. S. Masset, A. Morbidelli, A. Crida, and J. Ferreira. Disk Surface Density Transitions as Protoplanet Traps. *Astrophys. J.*, 642:478–487, May 2006.
- [41] F. S. Masset and J. C. B. Papaloizou. Runaway Migration and the Formation of Hot Jupiters. *Astrophys. J.*, 588:494–508, May 2003.
- [42] F. S. Masset and D. A. Velasco Romero. Dynamical friction on hot bodies in opaque, gaseous media. *Mon. Not. R. Astron. Soc.*, 465:3175–3184, March 2017.
- [43] K. Menou and J. Goodman. Low-Mass Protoplanet Migration in T Tauri α -Disks. *Astrophys. J.*, 606:520–531, May 2004.
- [44] N. Meyer-Vernet and B. Sicardy. On the physics of resonant disk-satellite interaction. *Icarus*, 69:157–175, January 1987.

- [45] A. F. Nelson and W. Benz. On the Early Evolution of Forming Jovian Planets. I. Initial Conditions, Systematics, and Qualitative Comparisons to Theory. *Astrophys. J.*, 589:556–577, May 2003.
- [46] R. P. Nelson. On the orbital evolution of low mass protoplanets in turbulent, magnetised disks. *Astron. Astrophys.*, 443:1067–1085, December 2005.
- [47] R. P. Nelson and J. C. B. Papaloizou. The interaction of a giant planet with a disc with MHD turbulence - II. The interaction of the planet with the disc. *Mon. Not. R. Astron. Soc.*, 339:993–1005, March 2003.
- [48] R. P. Nelson and J. C. B. Papaloizou. The interaction of giant planets with a disc with MHD turbulence - IV. Migration rates of embedded protoplanets. *Mon. Not. R. Astron. Soc.*, 350:849–864, May 2004.
- [49] R. P. Nelson, J. C. B. Papaloizou, F. S. Masset, and W. Kley. The migration and growth of protoplanets in protostellar discs. *Mon. Not. R. Astron. Soc.*, 318:18–36, October 2000.
- [50] G. I. Ogilvie and S. H. Lubow. Saturation of the Corotation Resonance in a Gaseous Disk. *Astrophys. J.*, 587:398–406, April 2003.
- [51] S.-J. Paardekooper. Dynamical corotation torques on low-mass planets. *Mon. Not. R. Astron. Soc.*, 444:2031–2042, November 2014.
- [52] S.-J. Paardekooper, C. Baruteau, and W. Kley. A torque formula for non-isothermal Type I planetary migration - II. Effects of diffusion. *Mon. Not. R. Astron. Soc.*, 410:293–303, January 2011.
- [53] S.-J. Paardekooper and G. Mellema. Halting type I planet migration in non-isothermal disks. *Astron. Astrophys.*, 459:L17–L20, November 2006.
- [54] S.-J. Paardekooper and J. C. B. Papaloizou. On disc protoplanet interactions in a non-barotropic disc with thermal diffusion. *Astron. Astrophys.*, 485:877–895, July 2008.
- [55] J. C. B. Papaloizou, R. P. Nelson, W. Kley, F. S. Masset, and P. Artymowicz. Disk-Planet Interactions During Planet Formation. In B. Reipurth, D. Jewitt, and K. Keil, editors, *Protostars and Planets V*, pages 655–668, 2007.
- [56] J. C. B. Papaloizou, R. P. Nelson, and M. D. Snellgrove. The interaction of giant planets with a disc with MHD turbulence - III. Flow morphology and conditions for gap formation in local and global simulations. *Mon. Not. R. Astron. Soc.*, 350:829–848, May 2004.
- [57] J. C. B. Papaloizou, R. P. Nelson, and M. D. Snellgrove. The interaction of giant planets with a disc with MHD turbulence - III. Flow morphology and conditions for gap formation in local and global simulations. *Mon. Not. R. Astron. Soc.*, 350:829–848, May 2004.
- [58] A. Pierens and J.-M. Huré. How does disk gravity really influence type-I migration? *Astron. Astrophys.*, 433:L37–L40, April 2005.
- [59] A. Pierens and R. P. Nelson. Constraints on resonant-trapping for two planets embedded in a protoplanetary disc. *Astron. Astrophys.*, 482:333–340, April 2008.

- [60] J. B. Pollack, O. Hubickyj, P. Bodenheimer, J. J. Lissauer, M. Podolak, and Y. Greenzweig. Formation of the Giant Planets by Concurrent Accretion of Solids and Gas. *Icarus*, 124:62–85, November 1996.
- [61] J. Szulagyi, F. Masset, E. Lega, A. Crida, A. Morbidelli, and T. Guillot. Circumplanetary disc or circumplanetary envelope? *Mon. Not. R. Astron. Soc.*, 460:2853–2861, August 2016.
- [62] H. Tanaka, T. Takeuchi, and W. R. Ward. Three-Dimensional Interaction between a Planet and an Isothermal Gaseous Disk. I. Corotation and Lindblad Torques and Planet Migration. *Astrophys. J.*, 565:1257–1274, February 2002.
- [63] C. E. J. M. L. J. Terquem. Stopping inward planetary migration by a toroidal magnetic field. *Mon. Not. R. Astron. Soc.*, 341:1157–1173, June 2003.
- [64] K. J. Walsh, A. Morbidelli, S. N. Raymond, D. P. O’Brien, and A. M. Mandell. A low mass for Mars from Jupiter’s early gas-driven migration. *Nature*, 475:206–209, July 2011.
- [65] W. R. Ward. Density waves in the solar nebula - Differential Lindblad torque. *Icarus*, 67:164–180, July 1986.
- [66] W. R. Ward. Protoplanet Migration by Nebula Tides. *Icarus*, 126:261–281, April 1997.
- [67] W. F. Winters, S. A. Balbus, and J. F. Hawley. Gap Formation by Planets in Turbulent Protostellar Disks. *Astrophys. J.*, 589:543–555, May 2003.

Evolución aproximada para un sistema compuesto por dos Hamiltonianos de Jaynes-Cummings acoplados

I. Ramos Prieto^{1,2}, A. Paredes¹, J. Récamier¹ y H. Moya Cessa²

¹ Instituto de Ciencias Físicas, UNAM, Apdo. Postal 48-3, Cuernavaca, Morelos 62251, México.

² Instituto Nacional de Astrofísica, Óptica y Electrónica, Calle Luis Enrique Erro No. 1, Santa María Tonanzintla, Puebla, 72840, México.

Abstract

En esta contribución presentamos la construcción aproximada del operador de evolución temporal para un sistema compuesto por dos Hamiltonianos tipo Jaynes-Cummings acoplados [1]. El hamiltoniano de Jaynes-Cummings consta de un campo unimodal y un átomo de dos niveles los cuales están acoplados por medio de interacciones que conservan el número total de excitaciones. Usando métodos de tipo algebraico expresamos el operador de evolución temporal como un producto de exponenciales y comparamos los resultados analíticos con resultados puramente numéricos para establecer la validez de nuestras aproximaciones.

1 Introducción

El modelo más simple para estudiar la interacción entre la materia y la luz cuantizada es el modelo de Jaynes-Cummings (JC), el cual tiene solución exacta debido a la aproximación de onda rotante [2] la cual consiste en descartar de la interacción aquellos términos que oscilan rápidamente. El modelo JC ha permitido la generación de estados no clásicos como son los estados de gato de Schrödinger [3], estados de número [4] y estados comprimidos [5]. Debido a que las propiedades algebraicas del modelo JC son similares a las que aparecen en la interacción ión-laser, se han propuesto diversas generalizaciones al modelo para estudiar efectos no lineales [6–8].

El enredamiento entre el estado de un átomo de dos niveles y el estado del campo en la cavidad se ha podido estudiar haciendo pasar el átomo

por la cavidad, ya que debido a la interacción entre ellos, el átomo, al salir de la cavidad lleva consigo información sobre el estado del campo y esta información puede extraerse por medio de mediciones condicionales sobre el átomo.

2 Teoría

Consideremos el acoplamiento de un sistema de dos niveles, modelado como un espín $1/2$, con un oscilador armónico cuántico. Esta es una situación típica en electrodinámica cuántica de cavidades (CQED) en donde un átomo se acopla con un modo de una cavidad [9]. El Hamiltoniano del sistema puede escribirse como

$$\hat{H} = \hat{H}_a + \hat{H}_c + \hat{H}_{ac}$$

en donde \hat{H}_a es al Hamiltoniano del átomo y \hat{H}_c es el hamiltoniano de la cavidad. El acoplamiento entre el átomo y la cavidad es $\hat{H}_{ac} = -\vec{D} \cdot \vec{E}_c$ en donde \vec{D} es el operador dipolar atómico y \vec{E}_c es el operador de campo eléctrico de la cavidad en la posición del átomo. El acoplamiento puede entonces escribirse como

$$\hat{H}_{ac} = -d[\epsilon_{\mathbf{a}}\hat{\sigma}_- + \epsilon_{\mathbf{a}}^*\hat{\sigma}_+] \cdot iE_0[\epsilon_{\mathbf{c}}\hat{a} - \epsilon_{\mathbf{c}}^*\hat{a}^\dagger].$$

Al desarrollar el producto escalar se obtienen cuatro términos. El proporcional a $\sigma_+ a^\dagger$ corresponde a una transición en donde el átomo pasa del estado $|g\rangle$ al estado $|e\rangle$ y además hay la creación de un fotón. El término $\sigma_- a$ corresponde a la transición $|e\rangle \rightarrow |g\rangle$ con la pérdida de un fotón. Cuando las frecuencias del campo y de la transición atómica son cercanas, estos procesos son altamente no resonantes y juegan un papel secundario en la evolución del sistema. Los otros dos términos corresponden a los procesos usuales de emisión y absorción de fotones. La aproximación de onda rotante consiste en conservar en la interacción solamente los términos resonantes con lo cual, el acoplamiento se reduce a

$$H_{ac} = -i\hbar \frac{\Omega_0}{2} [a\sigma_+ - a^\dagger\sigma_-], \quad (1)$$

en donde se ha introducido la frecuencia de Rabi del vacío

$$\Omega_0 = 2 \frac{dE_0 \epsilon_{\mathbf{a}}^* \cdot \epsilon_{\mathbf{c}}}{\hbar}.$$

Supondremos que $\epsilon_{\mathbf{a}}^* \cdot \epsilon_{\mathbf{c}}$ es real y positivo con lo cual, la frecuencia de Rabi también es real y positiva. La frecuencia de Rabi es una medida de la intensidad del acoplamiento átomo-campo. El Hamiltoniano del sistema puede escribirse como:

$$\hat{H} = \hbar \left(\omega_c \hat{a}^\dagger \hat{a} + \frac{\Omega}{2} \hat{\sigma}_z - i \frac{\Omega_0}{2} [\hat{a} \hat{\sigma}_+ - \hat{a}^\dagger \hat{\sigma}_-] \right) \quad (2)$$

En ausencia de interacción, los eigenestados del Hamiltoniano son el producto directo de estados del átomo y de estados del oscilador $|e, n\rangle$ y $|g, n\rangle$ y sus energías son $\hbar(\Omega/2 + n\omega_c)$ y $\hbar(-\Omega/2 + n\omega_c)$. Cuando el desentonamiento $\Delta = \Omega - \omega_c$ es pequeño, los estados $|g, n+1\rangle$ y $|e, n\rangle$ son casi degenerados y los niveles del sistema se pueden arreglar como una escalera de dobletes quedando el estado $|g, 0\rangle$ como el estado de más baja energía. Los diferentes dobletes corresponden a estados con $|g, n+1\rangle$ y $|e, n\rangle$ de forma que el número de excitaciones en el escalón es $M = n + 1$. El término de interacción solamente acopla los estados de un mismo escalón (los estados $|e, n\rangle$ y $|g, n+1\rangle$) por lo que el Hamiltoniano de Jaynes-Cummings conserva el número total de excitaciones M .

Podemos entonces rephrasar el problema y estudiar n sistemas de dos niveles, uno para cada valor de M .

El Hamiltoniano para el n -ésimo escalón tiene la forma:

$$\hat{H}_n = \hbar\omega_c \left(\hat{n} + \frac{1}{2} \right) \mathcal{I} + \hat{V}_n$$

en donde \mathcal{I} es la matriz identidad y

$$\hat{V}_n = \frac{\hbar}{2} \begin{pmatrix} \Delta & -i\Omega_n \\ i\Omega_n & -\Delta \end{pmatrix} = \frac{\hbar}{2} [\Delta \hat{\sigma}_z + \Omega_n \hat{\sigma}_y]$$

donde $\Omega_n = \sqrt{n+1} \Omega_0$ es la frecuencia de Rabi de n fotones.

El problema se reduce a diagonalizar una matriz de 2×2 . Los eigenvalores resultan

$$E_n^\pm = \hbar\omega_c \left(n + \frac{1}{2} \right) \pm \frac{\hbar}{2} \sqrt{\Delta^2 + \Omega_n^2} \quad (3)$$

con eigenvectores

$$\begin{aligned} |+, n\rangle &= \cos(\theta_n/2) |e, n\rangle + i \sin(\theta_n/2) |g, n+1\rangle \\ |-, n\rangle &= \sin(\theta_n/2) |e, n\rangle - i \sin(\theta_n/2) |g, n+1\rangle. \end{aligned} \quad (4)$$

donde $\tan \theta_n = \Omega_n / \Delta$.

Vemos pues que el problema tiene solución exacta gracias a la aproximación de onda rotante.

Consideremos ahora el Hamiltoniano

$$\hat{H} = \hbar \sum_{j=1}^2 \left[\omega_j \hat{a}_j^\dagger \hat{a}_j + \frac{\Omega_j}{2} \hat{\sigma}_z^{(j)} + g_j (\hat{a}_j \hat{\sigma}_+^{(j)} + \hat{a}_j^\dagger \hat{\sigma}_-^{(j)}) \right] + \hbar \lambda (\hat{a}_1 \hat{a}_2^\dagger + \hat{a}_1^\dagger \hat{a}_2), \quad (5)$$

que corresponde a dos Hamiltonianos tipo Jaynes-Cummings más un término de acoplamiento entre las cavidades. Tomamos como Hamiltoniano no perturbado

$$\hat{H}_0 = \hbar \left(\omega_1 \hat{n}_1 + \frac{\Omega_1}{2} \hat{\sigma}_z^{(1)} + \omega_2 \hat{n}_2 + \frac{\Omega_2}{2} \hat{\sigma}_z^{(2)} \right)$$

cuyo operador de evolución temporal es

$$\hat{U}_0 = e^{-i\omega_1 t \hat{n}_1} e^{-i\frac{\Omega_1}{2} t \hat{\sigma}_z^{(1)}} e^{-i\omega_2 t \hat{n}_2} e^{-i\frac{\Omega_2}{2} t \hat{\sigma}_z^{(2)}}$$

y el Hamiltoniano en la representación de interacción $\hat{H}_I(t)$ se obtiene de [10]:

$$\hat{H}_I(t) = \hat{U}_0^\dagger \hat{V} \hat{U}_0 \quad (6)$$

siendo

$$\hat{V} = \hbar \sum_{j=1}^2 g_j (\hat{a}_j \hat{\sigma}_+^j + \hat{a}_j^\dagger \hat{\sigma}_-^j) + \hbar \lambda (\hat{a}_1 \hat{a}_2^\dagger + \hat{a}_1^\dagger \hat{a}_2).$$

Al aplicar la transformación, obtenemos el Hamiltoniano en la representación de interacción

$$\begin{aligned} \hat{H}_I &= \hbar g_1 (\hat{\sigma}_+^{(1)} \hat{a}_1 e^{-i(\omega_1 - \Omega_1)t} + \hat{\sigma}_-^{(1)} \hat{a}_1^\dagger e^{i(\omega_1 - \Omega_1)t}) \\ &\quad + \hbar g_2 (\hat{\sigma}_+^{(2)} \hat{a}_2 e^{-i(\omega_2 - \Omega_2)t} + \hat{\sigma}_-^{(2)} \hat{a}_2^\dagger e^{i(\omega_2 - \Omega_2)t}) \\ &\quad + \hbar \lambda (\hat{a}_1^\dagger \hat{a}_2 e^{i(\omega_1 - \omega_2)t} + \hat{a}_1 \hat{a}_2^\dagger e^{-i(\omega_1 - \omega_2)t}) \end{aligned} \quad (7)$$

que separamos como $\hat{H}_I = \hat{V}_1 + \hat{V}_2$ con:

$$\hat{V}_1 = \hbar \lambda (\hat{a}_1^\dagger \hat{a}_2 e^{i(\omega_1 - \omega_2)t} + \hat{a}_1 \hat{a}_2^\dagger e^{-i(\omega_1 - \omega_2)t})$$

$$\begin{aligned} \hat{V}_2 &= \hbar g_1 (\hat{\sigma}_+^{(1)} \hat{a}_1 e^{-i(\omega_1 - \Omega_1)t} + \hat{\sigma}_-^{(1)} \hat{a}_1^\dagger e^{i(\omega_1 - \Omega_1)t}) \\ &\quad + \hbar g_2 (\hat{\sigma}_+^{(2)} \hat{a}_2 e^{-i(\omega_2 - \Omega_2)t} + \hat{\sigma}_-^{(2)} \hat{a}_2^\dagger e^{i(\omega_2 - \Omega_2)t}) \end{aligned}$$

El operador de evolución temporal en la representación de interacción satisface que $\hat{U}_I = \hat{U}_I^{(1)}\hat{U}_I^{(2)}$ en donde

$$i\hbar\partial_t\hat{U}_I^{(1)} = \hat{V}_1\hat{U}_I^{(1)}, \quad \hat{U}_I^{(1)}(t_0, t_0) = I \quad (8)$$

y

$$i\hbar\partial_t\hat{U}_I^{(2)} = \left[\hat{U}_I^{(1)\dagger}\hat{V}_2\hat{U}_I^{(1)} \right] \hat{U}_I^{(2)} \equiv \hat{H}_I^{(2)}\hat{U}_I^{(2)}, \quad \hat{U}_I^{(2)}(t_0, t_0) = I. \quad (9)$$

Para obtener el operador de evolución $\hat{U}_I^{(1)}$ definimos los operadores

$$\hat{J}_+ = \hat{a}_1\hat{a}_2^\dagger, \quad \hat{J}_- = \hat{a}_1^\dagger\hat{a}_2, \quad \hat{J}_z = \hat{a}_1^\dagger\hat{a}_1 - \hat{a}_2^\dagger\hat{a}_2$$

la interacción \hat{V}_1 está dada entonces por:

$$\hat{V}_1 = \hbar\lambda \left(\hat{J}_+ e^{-i(\omega_1 - \omega_2)t} + \hat{J}_- e^{i(\omega_1 - \omega_2)t} \right). \quad (10)$$

Las reglas de conmutación entre los operadores \hat{J}_i son:

$$[\hat{J}_-, \hat{J}_+] = \hat{J}_z, \quad [\hat{J}_-, \hat{J}_z] = -2\hat{J}_-, \quad [\hat{J}_+, \hat{J}_z] = 2\hat{J}_+$$

de donde vemos que el conjunto de operadores $\{\hat{J}_+, \hat{J}_-, \hat{J}_z\}$ es cerrado ante la operación de conmutación y el operador \hat{V}_1 es una combinación lineal de éstos.

Para escribir el operador de evolución temporal correspondiente usaremos el teorema de Wei-Norman [11] que establece que:

Si el Hamiltoniano del sistema puede escribirse como una combinación lineal de operadores \hat{X}_i

$$\hat{H} = \sum_{i=1}^N F_i(t)\hat{X}_i,$$

tal que para cualquier pareja \hat{X}_i, \hat{X}_j en \hat{H} el conmutador $[\hat{X}_i, \hat{X}_j] = \epsilon_{i,j}^k \hat{X}_k$ con $\epsilon_{i,j}^k$ una constante y \hat{X}_k un operador en \hat{H} , entonces el operador de evolución temporal \hat{U} puede escribirse como un producto de exponenciales

$$\hat{U} = \prod_{i=1}^N e^{\alpha_i(t)\hat{X}_i}$$

con $\alpha_i(t)$ funciones complejas, dependientes del tiempo por determinar. En el producto entran todos los elementos del algebra, no solamente los que aparezcan en \hat{H} .

Usando entonces el teorema de Wei-Norman escribimos el operador $\hat{U}_I^{(1)}$ como:

$$\hat{U}_I^{(1)} = e^{\gamma_1 \hat{J}_+} e^{\gamma_2 \hat{J}_-} e^{\gamma_3 \hat{J}_z}. \quad (11)$$

Para las funciones $\gamma_i(t)$ se obtiene el siguiente conjunto de ecuaciones diferenciales ordinarias al sustituir la Ec. 11 en la ecuación de Schrödinger:

$$\dot{\gamma}_1 = -i\lambda(e^{-i(\omega_1-\omega_2)t} - \gamma_1^2 e^{i(\omega_1-\omega_2)t}), \quad (12)$$

$$\dot{\gamma}_2 = -i\lambda(1 + 2\gamma_1\gamma_2)e^{i(\omega_1-\omega_2)t}, \quad (13)$$

$$\dot{\gamma}_3 = -i\lambda\gamma_1 e^{i(\omega_1-\omega_2)t}. \quad (14)$$

conocidos los coeficientes en el operador $\hat{U}_I^{(1)}$ podemos aplicar la transformación dada en la ecuación 9 para obtener el hamiltoniano en la nueva representación. Este está dado por:

$$\begin{aligned} \hat{H}_I^{(2)} = \hbar g_1 & \left[(1 + \gamma_1\gamma_2)e^{-\gamma_3 - i(\Omega_1 - \omega_1)t} \hat{a}_1^\dagger \hat{\sigma}_-^{(1)} + e^{\gamma_3 + i(\Omega_1 - \omega_1)t} \hat{a}_1 \hat{\sigma}_+^{(1)} \right] \\ & \hbar g_2 \left[(1 + \gamma_1\gamma_2)e^{-\gamma_3 + i(\Omega_2 - \omega_2)t} \hat{a}_2 \hat{\sigma}_+^{(2)} + e^{\gamma_3 - i(\Omega_2 - \omega_2)t} \hat{a}_2^\dagger \hat{\sigma}_-^{(2)} \right] \end{aligned} \quad (15)$$

en donde hemos despreciado términos que acoplan al átomo de la cavidad 1 con el campo de la cavidad 2 y viceversa. El Hamiltoniano $\hat{H}_I^{(2)}$ es una suma de dos Hamiltonianos tipo Jaynes-Cummings, uno para cada cavidad y átomo. Como ya mostramos antes, este problema tiene solución exacta (ver la ecuación 4). El operador de evolución temporal $\hat{U}_I^{(2)}$ puede entonces escribirse como el producto de operadores de evolución correspondientes a cada Hamiltoniano de Jaynes-Cummings, esto es

$$\hat{U}_I^{(2)} = \hat{U}_{JC}^{(1)} \hat{U}_{JC}^{(2)}. \quad (16)$$

Para encontrar la forma del operador de evolución temporal, introducimos los operadores:

$$\hat{b}_i = \frac{1}{\sqrt{\hat{M}_i}} \hat{a}_i \hat{\sigma}_+^{(i)}, \quad \hat{b}_i^\dagger = \hat{a}_i^\dagger \hat{\sigma}_-^{(i)} \frac{1}{\sqrt{\hat{M}_i}}$$

en donde el operador $\hat{M}_i = \hat{n}_i + \frac{1}{2}(1 + \hat{\sigma}_z^{(i)})$ nos da el número total de excitaciones en un escalón dado. Al actuar sobre los estados de la base, estos operadores dan:

$$\hat{b}_i |n_i, e_i\rangle = 0, \quad \hat{b}_i |n_i + 1, g_i\rangle = |n_i, e_i\rangle$$

$$\hat{b}_i^\dagger |n_i, e_i\rangle = |n_i + 1, g_i\rangle, \quad \hat{b}_i^\dagger |n_i + 1, g_i\rangle = 0$$

$$\hat{M}_i |n_i, e_i\rangle = (n_i + 1) |n_i, e_i\rangle, \quad \hat{M}_i |n_i + 1, g_i\rangle = (n_i + 1) |n_i + 1, g_i\rangle$$

y al aplicar \hat{b}_i^2 , $\hat{b}_i^{\dagger 2}$ sobre cualquier estado de la base, el resultado es cero. De las expresiones anteriores obtenemos los conmutadores:

$$[\hat{b}_i, \hat{b}_i^\dagger] = \hat{\sigma}_z^{(i)}, \quad [\hat{\sigma}_z^{(i)}, \hat{b}_i] = 2\hat{b}_i, \quad [\hat{\sigma}_z^{(i)}, \hat{b}_i^\dagger] = -2\hat{b}_i^\dagger. \quad (17)$$

El Hamiltoniano de interacción puede escribirse como

$$H_I^{(2)} = \hbar g_1 \sqrt{\hat{M}_1} [\phi_{11}(t) \hat{b}_1^\dagger + \phi_{12}(t) \hat{b}_1] + \hbar g_2 \sqrt{\hat{M}_2} [\phi_{21}(t) \hat{b}_2^\dagger + \phi_{22}(t) \hat{b}_2] \quad (18)$$

cuyo operador de evolución temporal es

$$\hat{U}_I^{(2)} = e^{\beta_z^{(1)} \hat{\sigma}_z^{(1)}} e^{\beta_+^{(1)} \hat{b}_1^\dagger} e^{\beta_-^{(1)} \hat{b}_1} e^{\beta_z^{(2)} \hat{\sigma}_z^{(2)}} e^{\beta_+^{(2)} \hat{b}_2^\dagger} e^{\beta_-^{(2)} \hat{b}_2} = \hat{U}_{JC}^{(1)} \hat{U}_{JC}^{(2)} \quad (19)$$

con funciones complejas dependientes del tiempo $\beta_i^{(j)}$ por determinar, las cuales cumplen con la condición inicial $\beta_i^{(j)}(t_0) = 0$. Nótese que para el estado $|0, g_i\rangle$, $M_i = 0$ y $\dot{\beta}_i^{(j)}(t_0) = 0$ y el operador $\hat{U}_I^{(2)}$ es el operador identidad.

El operador de evolución completo es entonces

$$\hat{U} = \hat{U}_0 \hat{U}_I^{(1)} \hat{U}_{JC}^{(1)} \hat{U}_{JC}^{(2)} \quad (20)$$

3 Resultados numéricos

Aplicaremos el operador temporal que hemos construido en la sección anterior al estado inicial $|\Psi(0)\rangle = |0, e_1\rangle \otimes |2, g_2\rangle$, esto es, la cavidad uno en el estado de vacío con un átomo en el estado excitado y la cavidad dos con dos fotones y un átomo en el estado base. El hamiltoniano de Jaynes-Cummings conserva el número de excitaciones en el sistema átomo-cavidad y el Hamiltoniano de interacción entre las cavidades conserva el número total de fotones (crea un fotón en una cavidad y aniquila un fotón en la otra), por lo tanto el número total de excitaciones del sistema es constante $M = M_1 + M_2$. Aplicando el operador de evolución temporal al estado inicial, obtenemos la función de onda al tiempo t como:

$$|\Psi(t)\rangle = \sum_{n_1=0}^3 \sum_{n_2=0}^3 \sum_{s_1=0}^1 \sum_{s_2=0}^1 \phi_{n_1, n_2}^{s_1, s_2}(t) |n_1, s_1\rangle \otimes |n_2, s_2\rangle \quad (21)$$

con la condición $n_1 + n_2 + s_1 + s_2 = 3$.

Para calcular la evolución del sistema fijamos los parámetros del Hamiltoniano de acuerdo a la referencia [12]. Las frecuencias de las cavidades $\omega_1/(2\pi) = 4$ GHz, $\omega_2/(2\pi) = 5$ GHz, las frecuencias de transición atómicas $\Omega_1 = 0.999\omega_1$ y $\Omega_2 = 0.999\omega_2$. El tiempo lo damos en unidades del periodo de la cavidad uno ($T_1 = .25 \times 10^{-9}$ s). Los acoplamientos átomo-campo y cavidad-cavidad los hicimos variar y se especifican en cada caso.

En la figura 1 se muestra la evolución temporal del número medio de fotones en la cavidad uno $\langle \hat{n}_1 \rangle$, el valor medio de $\hat{\sigma}_z^{(1)}$ y el valor medio del número de excitaciones en la cavidad uno $\langle \hat{M}_1 \rangle$ para acoplamientos átomo-campo $g_1 = 0.04\omega_1$ y $g_2 = 0.04\omega_2$ y para el acoplamiento cavidad-cavidad $\lambda = 0.001\omega_1$ de forma que el acoplamiento entre las cavidades es pequeño comparado con el acoplamiento átomo-campo en cada cavidad.

Vemos que el número total de excitaciones en la cavidad permanece constante, esto implica un acoplamiento débil entre las cavidades y un intercambio nulo de fotones entre ellas. Por otra parte, es claro que cuando el átomo decae al estado base el número de fotones en la cavidad aumenta en la unidad (tomamos la frecuencia del campo prácticamente igual a la frecuencia de transición atómica) y viceversa. Este comportamiento es el esperado para un Hamiltoniano de Jaynes-Cummings.

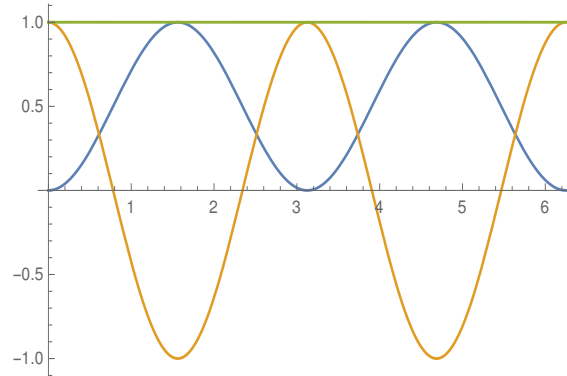


Figure 1: Valor medio del número de fotones en la cavidad uno $\langle \hat{n}_1 \rangle$ (azul) de $\langle \hat{\sigma}_z^{(1)} \rangle$ (anaranjado) y de $\langle \hat{M}_1 \rangle$ (verde). Parámetros $g_1 = 0.04\omega_1$, $g_2 = 0.04\omega_2$, $\lambda = 0.001\omega_1$.

En la figura 2 mostramos la evolución temporal de $\langle \hat{n}_1 \rangle$ y $\langle \hat{n}_2 \rangle$ con parámetros

$g_1 = 0.001\omega_1$, $g_2 = 0.001\omega_2$ y $\lambda = 0.25\omega_1$. En este caso el acoplamiento entre las cavidades es mucho más intenso que el acoplamiento átomo-cavidad, de forma que el intercambio de excitaciones entre el átomo y la cavidad es pequeño. Hay un importante intercambio de fotones entre las dos cavidades como puede verse claramente en la figura. Con el objeto de comprobar la

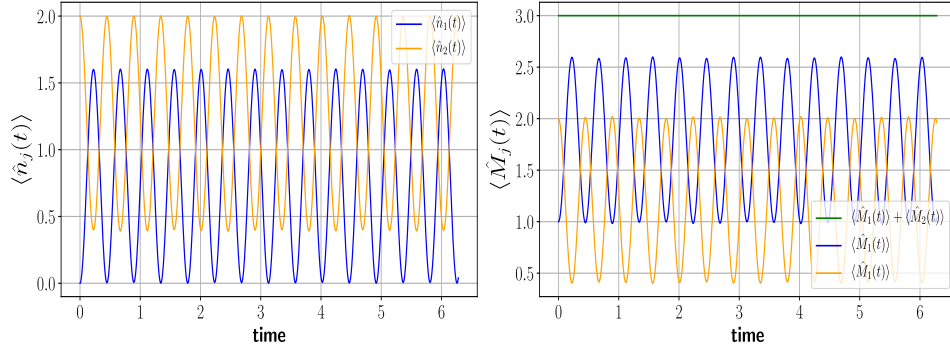


Figure 2: Panel izquierdo: Valor medio del número de fotones $\langle \hat{n}_1 \rangle$ (azul) y $\langle \hat{n}_2 \rangle$ (anaranjado). Panel derecho: Número total de excitaciones $\langle \hat{M}_1 \rangle$ (azul), $\langle \hat{M}_2 \rangle$ (anaranjado) y $\langle \hat{M}_1 + \hat{M}_2 \rangle$ (verde). Parámetros $g_1 = 0.001\omega_1$, $g_2 = 0.001\omega_2$, $\lambda = 0.25\omega_1$.

calidad de nuestras aproximaciones, hicimos también un cálculo puramente numérico de las cantidades reportadas en las figuras tomando una base del tamaño necesario para alcanzar convergencia [13]. Los resultados analíticos ajustan muy bien a los resultados numéricos lo cual indica que si la relación entre los parámetros es tal que uno puede considerarse pequeño comparado con el otro, entonces las aproximaciones utilizadas son adecuadas.

Finalmente consideremos el caso en el que los parámetros de acoplamiento átomo-campo y cavidad-cavidad son del mismo orden de magnitud. Como ejemplo tomemos $g_1 = 0.04\omega_1$, $g_2 = 0.04\omega_2$ y $\lambda = 0.08\omega_1$. Los resultados analíticos y numéricos se muestran en la figura 3. En el panel izquierdo se muestran $\langle \hat{n}_1 \rangle$ y $\langle \hat{n}_2 \rangle$ y en el derecho $\langle \hat{\sigma}_z^{(1)} \rangle$ y $\langle \hat{\sigma}_z^{(2)} \rangle$. El comportamiento en la evolución temporal del número de fotones en cada cavidad difiere apreciablemente del encontrado en los casos anteriores. En este caso, el acoplamiento entre los átomos y la cavidad correspondiente es grande, así también el acoplamiento entre las cavidades. En el panel derecho vemos que los promedios analíticos de la evolución atómica son funciones oscilatorias

con amplitud casi constante. La solución numérica tiene un comportamiento cualitativo similar pero se puede observar un desfase entre los dos cálculos conforme el tiempo avanza. En el panel izquierdo vemos que el comportamiento cualitativo entre el cálculo numérico y el cálculo analítico para los valores medios de fotones en cada cavidad es también similar, especialmente a tiempos cortos.

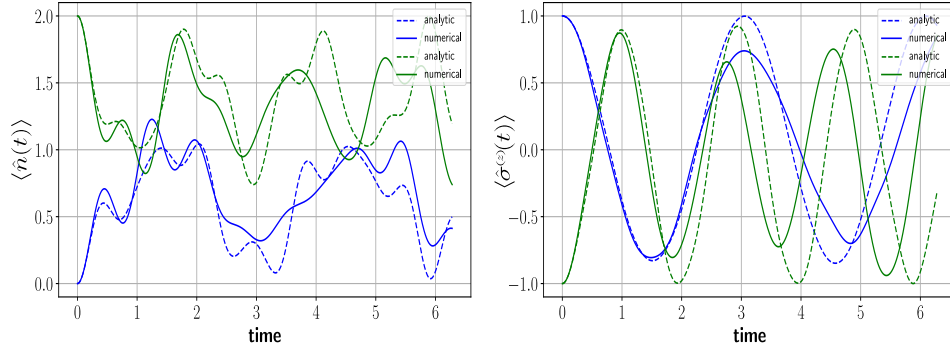


Figure 3: Panel izquierdo: Valor medio del número de fotones $\langle \hat{n}_1 \rangle$ (azul) y $\langle \hat{n}_2 \rangle$ (verde). Panel derecho: Valor medio del estado atómico $\langle \hat{\sigma}_z^{(1)} \rangle$ (azul) y $\langle \hat{\sigma}_z^{(2)} \rangle$ (verde). Parámetros $g_1 = 0.04\omega_1$, $g_2 = 0.04\omega_2$, $\lambda = 0.08\omega_1$.

4 Conclusiones

En este trabajo hemos presentado una forma aproximada para construir el operador de evolución temporal correspondiente a un sistema compuesto por dos Hamiltonianos acoplados tipo Jaynes-Cummings. El acoplamiento genera una excitación en una cavidad y aniquila una excitación en la otra cavidad. Cuando solamente se tiene el Hamiltoniano de Jaynes-Cummings el problema tiene solución exacta. Cuando se tienen las cavidades acopladas sin átomos en éstas el problema también tiene solución exacta. Sin embargo, cuando se tiene el acoplamiento entre las cavidades y también se tienen átomos en cada cavidad el problema no tiene solución exacta. Para resolver el problema pasamos a la representación de interacción y aproximamos el Hamiltoniano en dicha representación conservando los operadores que cierran un algebra de Lie. De esta manera es posible obtener un operador de evolución temporal exacto para una interacción aproximada. Esto

es importante ya que el operador de evolución conserva sus propiedades de unitariedad. Una vez que escribimos el operador de evolución del sistema completo en forma de un producto de exponenciales es posible hacer evolucionar el sistema y obtener la matriz densidad a un tiempo t cualquiera con lo cual se tienen todas las herramientas necesarias para calcular cualquier propiedad que nos interese como el valor medio del número de fotones, la función de Husimi, etc.

Para asegurarnos de la validez de nuestras aproximaciones hicimos un cálculo puramente numérico de las cantidades que evaluamos analíticamente y encontramos un excelente acuerdo entre ambos cálculos especialmente en aquellos casos en donde los parámetros de acoplamiento cavidad-cavidad y átomo-cavidad difieren significativamente. Para el caso en que los parámetros son del mismo orden encontramos una buena concordancia a nivel cualitativo.

Agradecimientos Agradecemos a la DGAPA UNAM por su apoyo através del proyecto PAPIIT IN111119.

References

- [1] I Ramos-Prieto, A. Paredes, J. Récamier and H. Moya-Cessa, 2019 Phys. Scr. in press <https://doi.org/10.1088/1402-4896/ab538b>
- [2] E. T. Jaynes and F. W. Cummings, Proc. IEEE **51**, 89 (1963).
- [3] K. Vogel, V. M. Akulin and W. P. Schleich, Phys. Rev. Lett. **71**, 1816 (1993).
- [4] H. Moya-Cessa, P. L. Knight and A. Rosenhouse-Dantsker, Phys. Rev. A **50**, 1814-1821 (1994).
- [5] J. R. Kuklinski, J. L. Madajczyk, Phys. Rev. A **37**, 3175-3178 (1988).
- [6] S. Cordero and J. Récamier, J. Phys. A:Math. Theor. **45**, 385303 (2012).
- [7] O. de los Santos-Sánchez, J. Récamier, J. Phys. B: At. Mol. Opt. Phys. **45**, 015502 (2012).
- [8] R. L. Matos Filho and W. Vogel, Phys. Rev. Lett. **76**, 608 (1996).
- [9] *Exploring the Quantum: Atoms, Cavities, and Photons*, Serge Haroche and Jean-Michel Raimond, Oxford Graduate Texts, (2006).

- [10] *Modern Quantum Mechanics*, Sakurai, Napolitano, Addison Wesley (2011).
- [11] J Wei and E Norman, *Journal of Mathematical Physics* **4**, 575 (1963).
- [12] S. Felicetti, M. Sanz, L. Lamata, G. Romero, G. Johansson, P. Delsing and E. Solano, *Phys. Rev. Lett.* **113**, 093602 (2014).
- [13] J. R. Johansson, P. D. Nation, and F. Nori, "QuTiP2: A Python framework for the dynamics of open quantum systems", *Comp. Phys. Comm.* **184**, 1234 (2013).

A brief introduction to the NEGF method for electron transport at the nanoscale*

Thomas Stegmann
Instituto de Ciencias Físicas
Universidad Nacional Autónoma de México

November 19, 2019

1 Introduction

On the long way of transport theory in solid state materials, one of the first milestones was the Drude-Sommerfeld model [2, 3], which allows to understand electric and thermal conduction in bulk metals at room temperature. Its main ideas are that the electrons follow the Fermi-Dirac distribution and only electrons close to the Fermi energy contribute to the conduction. These conduction electrons are scattered randomly after a given length. This gives rise to Ohm's law in metals. Further significant steps were Bloch's theorem [4] and the band structure theory [5–8], which allow to understand the difference between metals, semiconductors and insulators.

However, when the system size is reduced to the nanometer scale, several characteristic lengths have the same order of magnitude:

- Mean free path ℓ_m : The average distance after which an electron is scattered and its momentum is randomized.
- Phase coherence length ℓ_ϕ : The average distance after which the electron phase is randomized.
- Fermi wavelength λ_F : Wavelength of the electrons at the Fermi energy. The Fermi wavelength gives the length scale at which quantum effects emerge.
- Localization length λ : Average spatial extent of the exponentially decaying electron eigenstates in a disordered quantum system.
- System size L

*These lecture notes are based on Ref. [1].

This gives rise to novel conduction phenomena. An electron can propagate *ballistically* through a nanosystem without any scattering event, which results in a length-independent resistance caused solely by the contacts [9]. Ballistic transport is found, for example, in carbon nanotubes [10, 11] and graphene [12, 13]. Moreover, when not only the momentum of the electrons is preserved but also its phase information, *interference effects* may influence the transport drastically, as it can be seen for example in Aharonov-Bohm experiments [14, 15] and in the coherent electron focusing [16–19]. In contrast, in disordered nanosystems *Anderson localization* [20–23] is found, which manifests in an exponential increase of the resistance with the system size. Therefore, in order to describe the transport processes on the nanometer scale correctly, a completely different approach is necessary, which starts with a microscopic quantum description of the system.

In these lecture notes, we give a short but self-contained introduction into quantum transport using the non-equilibrium Green’s function (NEGF) approach. Instead of applying many-body perturbation theory [24, 25], we follow the textbooks by Datta [9, 26, 27] and motivate physically the required equations.

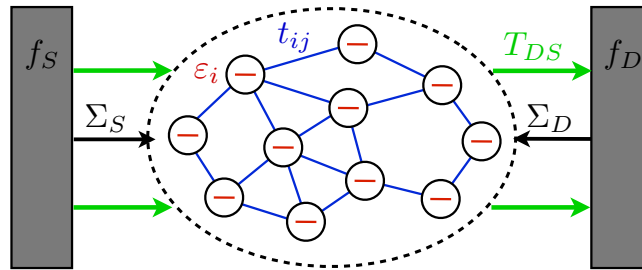


Figure 1: The transport of electrons from the source S to the drain D through a nanosystem is studied. The nanosystem is described by the Hamiltonian (1) and can be visualized by a network of sites ε_i and bonds t_{ij} , see the region enclosed by the dashed ellipse.

2 Definitions and concepts

2.1 Hamiltonian of the nanosystem

The nanosystem is described by a tight-binding Hamiltonian

$$H = \sum_i \varepsilon_i |i\rangle \langle i| + \sum_{j<i} (t_{ij} |i\rangle \langle j| + \text{H.c.}), \quad (1)$$

where H.c. means Hermitian conjugation. This Hamiltonian can be visualized by a network of sites and bonds, see the region enclosed by the dashed ellipse in Figure 1. The onsite energy ε_i represents the potential energy, which is necessary to occupy the state

$|i\rangle$ (the i th site) by an electron. The coupling matrix element t_{ij} is proportional to the transition rate of electrons from state $|i\rangle$ to state $|j\rangle$ (from the i th site to the j th site). For a finite system of N sites the Hamiltonian (1) can be represented as a $N \times N$ matrix with diagonal matrix elements ϵ_i and off-diagonal elements t_{ij} .

2.2 The spectral function A

The spectral function is defined as

$$A(E) \equiv 2\pi\delta(E - H), \quad (2)$$

where E is the electron energy and H the tight-binding Hamiltonian¹. In real space the expansion in terms of eigenfunctions $\psi_k(\mathbf{r})$ and eigenenergies ϵ_k of the Hamiltonian H

$$A(\mathbf{r}, \mathbf{r}', E) = 2\pi \sum_k \psi_k(\mathbf{r})\delta(E - \epsilon_k)\psi_k^*(\mathbf{r}') \quad (3)$$

shows that the diagonal elements of the spectral function (apart from the factor 2π) give the local density of states $D(\mathbf{r}, E)$ (LDOS). The total density of states $D(E)$ (DOS) is then obtained by integration over space, or more generally, by the trace of (2)

$$D(E) = \frac{1}{2\pi} \text{Tr}(A(E)) = \sum_k \delta(E - \epsilon_k). \quad (4)$$

2.3 The Green's functions G and G^+

The δ distribution can be represented by

$$2\pi\delta(E - \epsilon_k) = \frac{2\nu}{(E - \epsilon_k)^2 + \nu^2} = i \left[\frac{1}{E - \epsilon_k + i\nu} - \frac{1}{E - \epsilon_k - i\nu} \right], \quad (5)$$

where ν is an infinitesimal positive number. Using this representation of the δ distribution, the spectral function can be written as

$$A(E) = 2\pi\delta(E - H) = i \left[\underbrace{(E - H + i\nu)^{-1}}_{\substack{\text{retarded} \\ \text{Green's function } G}} - \underbrace{(E - H - i\nu)^{-1}}_{\substack{\text{advanced} \\ \text{Green's function } G^+}} \right] = -2 \text{Im}(G), \quad (6)$$

where we have introduced the (retarded) Green's function G and the advanced Green's function G^+ . As the DOS has to be positive, we also learn from (6) that the diagonal matrix elements of the Green's function fulfill $\text{Im}(G_{ii}) < 0$.

¹We omit unit matrices, which match scalars (energy E) to matrices (Hamiltonian H).

2.4 The correlation function G^n

We define the correlation function² as

$$G^n \equiv 2\pi |\psi\rangle \langle\psi|. \quad (7)$$

Its matrix elements in real space

$$G^n(\mathbf{r}, \mathbf{r}') = 2\pi \psi(\mathbf{r})\psi^*(\mathbf{r}'), \quad (8)$$

give the correlation of the state $|\psi\rangle$ between the places \mathbf{r} and \mathbf{r}' . In particular, its diagonal matrix elements ($\mathbf{r} = \mathbf{r}'$) give the electron density (apart from a factor 2π). In equilibrium the electron density is determined by the occupation of the density of states according to the Fermi distribution $f(E - \mu)$ with chemical potential μ . Therefore, we conclude that in equilibrium the correlation function is related to the spectral function by

$$G_{\text{eq}}^n(E) = A(E)f(E - \mu). \quad (9)$$

3 Open quantum systems: Σ and Σ^{in}

In order to study transport of electrons through the nanosystem, we have to connect it to a source S and a drain D reservoir, see the dark gray rectangles in Figure 1. The source and drain are in equilibrium and characterized by Fermi distributions $f_{S/D} \equiv f(E - \mu_{S/D})$ with chemical potentials $\mu_{S/D}$. Their difference drives the system out of equilibrium and causes the current flow.

The isolated reservoirs are described by the Hamiltonians $H_{S/D}$, which fulfill the Schrödinger equations

$$(E - H_S) |\Phi_S\rangle = 0, \quad (10a)$$

$$(E - H_D) |\Phi_D\rangle = 0. \quad (10b)$$

These equations can be rewritten in the form

$$(E - H_S + i\nu) |\Phi_S\rangle = |Q_S\rangle, \quad (11a)$$

$$(E - H_D + i\nu) |\Phi_D\rangle = |Q_D\rangle, \quad (11b)$$

where ν is an infinitesimal positive number. The term $i\nu |\Phi_{S/D}\rangle$ represents the extraction of electrons from the contact, whereas $|Q_{S/D}\rangle$ represents the reinjection of electrons from external sources. Extraction and reinjection are necessary to maintain the reservoirs in equilibrium. The Schrödinger equation is mathematically unchanged, if we identify $i\nu |\Phi_{S/D}\rangle = |Q_{S/D}\rangle$. However, the transition from (10) to (11) is not only a formal modification of the Schrödinger equation but a change in the point of view. In the latter, E

²Instead of the correlation function G^n , the lesser Green's function $G^<$ is also commonly used in the literature, see e.g. [28]. These functions are connected by $G^< = iG^n$.

is no longer an eigenenergy of the Hamiltonian but an independent variable, which gives the energy of excitations $|Q_{S/D}\rangle$ from external sources. Whereas in (10) the $|\Phi_{S/D}\rangle$ are non-zero only for the eigenenergies, in (11) the $|\Phi_{S/D}\rangle$ are non-zero for any energy and represent the response of the reservoirs to external excitations.

Now, what happens when the reservoirs are connected to the nanosystem by coupling matrices $\tau_{S/D}$? The states $|\Phi_{S/D}\rangle$ in the reservoirs spill over and excite states $|\psi\rangle$ in the nanosystem, which in return also excite states $|\chi_{S/D}\rangle$ in the reservoirs. The Schrödinger equation of the coupled system reads

$$\begin{pmatrix} E - H_S + i\nu & -\tau_S^+ & 0 \\ -\tau_S & E - H & -\tau_D \\ 0 & -\tau_D^+ & E - H_D + i\nu \end{pmatrix} \begin{pmatrix} \Phi_S + \chi_S \\ \psi \\ \Phi_D + \chi_D \end{pmatrix} = \begin{pmatrix} Q_S \\ 0 \\ Q_D \end{pmatrix}. \quad (12)$$

As it can be assumed that the reinjection $|Q_{S/D}\rangle$ is unchanged by the coupling, the first and the last row of (12) lead with (11) to

$$|\chi_S\rangle = G_S \tau_S^+ |\psi\rangle, \quad (13a)$$

$$|\chi_D\rangle = G_D \tau_D^+ |\psi\rangle, \quad (13b)$$

where

$$G_S \equiv (E - H_S + i\nu)^{-1}, \quad (14a)$$

$$G_D \equiv (E - H_D + i\nu)^{-1} \quad (14b)$$

are the Green's functions of the reservoirs. By means of (13) the middle row of (12) reads

$$(E - H - \Sigma_S - \Sigma_D) |\psi\rangle = |Q\rangle, \quad (15)$$

where we defined the so-called **self-energies**

$$\Sigma_S \equiv \tau_S G_S \tau_S^+, \quad (16a)$$

$$\Sigma_D \equiv \tau_D G_D \tau_D^+ \quad (16b)$$

and the total excitation of the nanosystem

$$|Q\rangle \equiv \tau_S |\Phi_S\rangle + \tau_D |\Phi_D\rangle. \quad (17)$$

Finally, we can write for its states

$$|\psi\rangle = G |Q\rangle, \quad (18)$$

where we defined the **Green's function of the nanosystem**

$$G \equiv (E - H - \Sigma_S - \Sigma_D)^{-1}. \quad (19)$$

Therefore, the Schrödinger equation of the coupled system has been transformed to a single equation for the nanosystem, which is “open” to the environment by self-energies. This approach simplifies the problem drastically because the dimension of the Hilbert space of the nanosystem is much smaller than the dimension of the Hilbert space of the coupled system. The self-energies represent a non-Hermitian modification of the Hamiltonian, which shift its eigenenergies from the real axis into the complex plane. The imaginary part of the eigenenergies is inversely proportional to the lifetime of the states in the nanosystem and causes an energy broadening, see Figure 2.

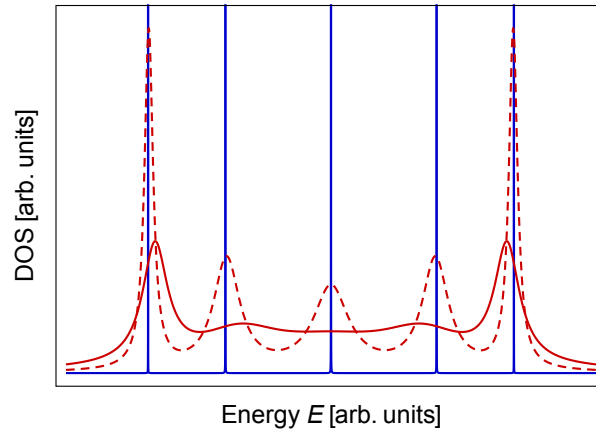


Figure 2: Energy resolved density of states (DOS) of a typical nanosystem. The DOS of the isolated system (blue curve) shows peaks at the discrete eigenenergies of the Hamiltonian. These peaks are broadened, if the system is opened increasingly to reservoirs (dashed to solid red curve).

With the definition of the **broadening matrix**

$$\Gamma \equiv i(\Sigma - \Sigma^+) \quad (20)$$

we obtain³ for the spectral function of the nanosystem

$$A \equiv i(G - G^+) = G(\Gamma_S + \Gamma_D)G^+ = A_1 + A_2. \quad (21)$$

These spectral functions $A_{1/2} \equiv G\Gamma_{S/D}G^+$ give the density of states for electrons originating S/D and should not be confused with the spectral functions $A_{S/D} = i(G_{S/D} - G_{S/D}^+)$, which give the density of states in the reservoirs. In order to calculate the correlation

³To obtain (21), we use $i((G^{-1})^+ - G^{-1}) = \Gamma_S + \Gamma_D$ and by multiplication with G from the left and G^+ from the right, we arrive at $i(G - G^+) = G(\Gamma_S + \Gamma_D)G^+$.

function, we study the projector

$$\begin{aligned}
|\psi\rangle\langle\psi| &\stackrel{(18)}{=} G|Q\rangle\langle Q|G^+ \\
&\stackrel{(17)}{=} G\tau_S|\Phi_S\rangle\langle\Phi_S|\tau_S^+G^+ + G\tau_D|\Phi_D\rangle\langle\Phi_D|\tau_D^+G^+ \\
&\quad + \underbrace{G\tau_S|\Phi_S\rangle\langle\Phi_D|\tau_D^+G^+ + G\tau_D|\Phi_D\rangle\langle\Phi_S|\tau_S^+G^+}_{=0, \text{ because no direct coupling between the reservoirs}}.
\end{aligned} \tag{22}$$

Applying (7) and using (9) for the reservoirs, in which equilibrium is assumed, we obtain for the non-equilibrium **correlation function of the nanosystem**

$$\begin{aligned}
G^m &= G\underbrace{\tau_S A_S \tau_S^+}_{\Gamma_S} G^+ f(E - \mu_S) + G\underbrace{\tau_D A_D \tau_D^+}_{\Gamma_D} G^+ f(E - \mu_D) \\
&= A_1 f(E - \mu_S) + A_2 f(E - \mu_D)
\end{aligned} \tag{23a}$$

$$= G\Sigma^{\text{in}}G^+, \tag{23b}$$

where we defined the **inscattering function**

$$\Sigma^{\text{in}} \equiv \Sigma_S^{\text{in}} + \Sigma_D^{\text{in}} = \Gamma_S f_S + \Gamma_D f_D. \tag{24}$$

As it could be expected for non-interacting electrons, the correlation function of the nanosystem, which gives the electron density, is the sum of the spectral functions occupied by the Fermi distributions of the corresponding contacts.

The equations for the Green's function (19) and for the correlation function (23b) are two essential results of the non-equilibrium Green's function approach. We have motivated them physically, but arrived essentially at the same result as in Keldysh's seminal paper [24, eqs. (75)–(77)], where many-body perturbation theory is applied. In fact the strength of the non-equilibrium Green's function approach is more profound because it allows to include arbitrary interactions in the system by suitable self-energies and inscattering functions, see e.g. [9, 28] for details. These interactions can formally be considered as additional virtual reservoirs attached to the system.

4 Current equations

In order to define the current operator, we start with the isolated nanosystem. By means of the time-dependent Schrödinger equation, we obtain for the time evolution of the projector

$$\frac{d}{dt} |\psi\rangle\langle\psi| + \frac{i}{\hbar} [H, |\psi\rangle\langle\psi|] = 0. \tag{25}$$

In the same way as the continuity equation of quantum mechanics,⁴ this equation reflects the conservation of electron density. However, when the nanosystem is connected to

⁴Indeed, when we go to local space by multiplying (25) with $\langle\mathbf{r}|$ from left and $|\mathbf{r}\rangle$ from right, we obtain the continuity equation in its common form.

reservoirs, electrons can enter and leave it and thus, the density in the system is not conserved. In steady state, the first term of (25) vanishes and the remaining commutator tells us the rate, at which electrons are lost in the system. Therefore, applying (7), the current operator can be defined as

$$I^{\text{op}} \equiv \frac{ie}{\hbar} [H, G^n], \quad (26)$$

where e is the electron charge. Its diagonal elements

$$I_{ii}^{\text{op}} = \frac{ie}{\hbar} \sum_j (t_{ij} G_{ji}^n - t_{ji} G_{ij}^n) \quad (27)$$

give the total current flowing to the i th site. Hence, the individual terms in this sum can be identified as the local current flowing between the i th and j th site [29, 30],

$$I_{ij}(E) = \frac{ie}{\hbar} (t_{ij} G_{ji}^n - t_{ji} G_{ij}^n) = \frac{2e}{\hbar} \text{Im} (t_{ij}^* G_{ij}^n). \quad (28)$$

The total flow of electrons with energy E through the dashed ellipse in Figure 1 is then given by⁵

$$I(E) \equiv \text{Tr} (I^{\text{op}}) = \frac{e}{\hbar} \text{Tr} (\Sigma^{\text{in}} A - \Gamma G^n). \quad (30)$$

As the number of electrons is conserved, the inflow of electrons equals the outflow and hence, the total flow is exactly zero. Taking into account the invariance of the trace under cyclic permutations, this can be seen already by means of the definition of the current operator (26). However, separating the inscattering function and the broadening function into the individual contributions of each reservoir, we obtain for the non-zero current at the drain

$$I_D(E) = \frac{e}{\hbar} \text{Tr} (\Sigma_D^{\text{in}} A - \Gamma_D G^n). \quad (31)$$

This current equation is another key result of the non-equilibrium Green's function approach. For non-interacting electrons it can be further simplified by means of (23a) and (24). Integrating over energy, we arrive finally at the famous **Landauer formula** [31, 32] for the total current through the system⁶

$$I_D = \frac{e}{\hbar} \int dE T_{DS}(E) (f(E - \mu_S) - f(E - \mu_D)), \quad (32)$$

⁵The current operator can be rewritten by means of

$$[H, G^n] \stackrel{(23b)}{=} HG\Sigma^{\text{in}}G^+ - G\Sigma^{\text{in}}G^+H \quad (29a)$$

$$\stackrel{(19)}{=} G\Sigma^{\text{in}} - \Sigma^{\text{in}}G^+ + \underbrace{G\Sigma^{\text{in}}G^+}_{G^n} \Sigma^+ - \Sigma \underbrace{G\Sigma^{\text{in}}G^+}_{G^n}. \quad (29b)$$

Using (20), (21) and the invariance of the trace under cyclic permutations, we obtain (30).

⁶If not stated otherwise, integration boundaries extend from $-\infty$ to $+\infty$.

where we defined the **transmission function**

$$T_{DS} \equiv \text{Tr} (\Gamma_D G \Gamma_S G^+) . \quad (33)$$

The transmission function T_{DS} gives the probability that an electron injected by the source will transmit to the drain. For an isolated system⁷, the transmission is perfect ($T_{DS} = 1$) at its eigenenergies, whereas it vanishes elsewhere, see Figure 3. These transmission peaks are broadened, if the system is increasingly opened to reservoirs. Note that for sufficiently strong coupling to the reservoirs, the transmission in the band center is perfect and nearly constant.

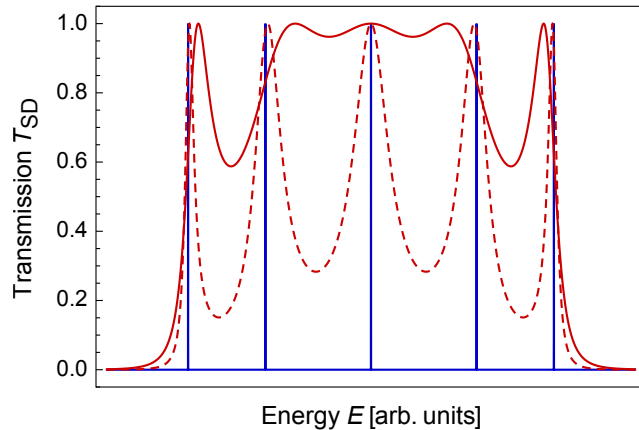


Figure 3: Transmission through the same nanosystem as in Figure 2. In the same way as the DOS, the transmission of the isolated system is perfect ($T_{DS} = 1$) at its eigenenergies, whereas it vanishes elsewhere (blue curve). These transmission peaks are broadened, if the system is opened increasingly to reservoirs (dashed to solid red curve).

5 Further reading

In these lecture notes, we have motivated physically the essential equations required to study electron transport at the nanoscale. However, to get a deeper understanding, it is indispensable to study the existing literature in detail. From the numerous textbooks, we would like to recommend some few. Maybe the best starting point are the textbooks by Datta [9, 26, 27], who gives a didactically brilliant introduction into the physics of nanosystems. A mathematically more rigorous and yet readable introduction can be found

⁷In order to define the transmission function, the isolated system is still infinitesimally weakly coupled to reservoirs.

in the book by Di Ventra [28]. The usage of Green's function in quantum mechanics is nicely explained in the book by Economou [33]. To get started the reader may also find useful the books by Heikkilä [34], Ihn [35] and Ferry, Goodnick and Bird [36]. We can also recommend the book on molecular electronics by Cuevas and Scheer [37].

References

- [1] T. Stegmann. *Quantum transport in nanostructures: From the effects of decoherence on localization to magnetotransport in two-dimensional electron systems*. PhD thesis, University Duisburg-Essen, 2014.
- [2] P. Drude. Zur Elektronentheorie der Metalle. *Ann. Phys.*, 306:566, 1900.
- [3] A Sommerfeld and H. Bethe. Elektronentheorie der Metalle. In *Handbuch der Physik*, volume 24, pages 333–622. Springer, 1933.
- [4] F. Bloch. Über die Quantenmechanik der Elektronen in Kristallgittern. *Z. Phys.*, 52:555, 1929.
- [5] R. Peierls. Zur Theorie der elektrischen und thermischen Leitfähigkeit von Metallen. *Ann. Phys.*, 4:121, 1930.
- [6] L. Brillouin. Les électrons libres dans les métaux et le role des réflexions de Bragg. *J. Phys. Radium*, 1:377, 1930.
- [7] A. H. Wilson. The theory of electronic semi-conductors. *Proc. R. Soc. Lond. A*, 133:458, 1931.
- [8] A.H. Wilson. The theory of electronic semi-conductors. ii. *Proc. R. Soc. Lond. A*, 134:277, 1931.
- [9] S. Datta. *Electronic Transport in Mesoscopic Systems*. Cambridge University Press, 1997.
- [10] S. Frank, P. Poncharal, Z. L. Wang, and Walt A. de Heer. Carbon nanotube quantum resistors. *Science*, 280:1744, 1998.
- [11] A. Bachtold, M. S. Fuhrer, S. Plyasunov, M. Forero, E. H. Anderson, A. Zettl, and P. L. McEuen. Scanned probe microscopy of electronic transport in carbon nanotubes. *Phys. Rev. Lett.*, 84:6082, 2000.
- [12] X. Du, I. Skachko, A. Barker, and E. Y. Andrei. Approaching ballistic transport in suspended graphene. *Nat. Nano.*, 3:491, 2008.
- [13] K. I. Bolotin, K. J. Sikes, J. Hone, H. L. Stormer, and P. Kim. Temperature-dependent transport in suspended graphene. *Phys. Rev. Lett.*, 101:096802, 2008.

- [14] A. Bachtold, C. Strunk, J.-P. Salvetat, J.-M. Bonard, L. Forro, T. Nussbaumer, and C. Schonenberger. Aharonov-bohm oscillations in carbon nanotubes. *Nature*, 397:673, 1999.
- [15] B. Grbić, R. Leturcq, T. Ihn, K. Ensslin, D. Reuter, and A. D. Wieck. Aharonov–bohm oscillations in p-type gaas quantum rings. *Physica E*, 40:1273, 2008.
- [16] H. van Houten, C. W. J. Beenakker, J. G. Williamson, M. E. I. Broekaart, P. H. M. van Loosdrecht, B. J. van Wees, J. E. Mooij, C. T. Foxon, and J. J. Harris. Coherent electron focusing with quantum point contacts in a two-dimensional electron gas. *Phys. Rev. B*, 39:8556, 1989.
- [17] K. E. Aidala, R. E. Parrott, T. Kramer, E. J. Heller, R. M. Westervelt, M. P. Hanson, and A. C. Gossard. Imaging magnetic focusing of coherent electron waves. *Nat. Phys.*, 3:464, 2007.
- [18] T. Stegmann, D. E. Wolf, and A. Lorke. Magnetotransport along a boundary: From coherent electron focusing to edge channel transport. *New J. Phys.*, 15:113047, 2013.
- [19] T. Stegmann and A. Lorke. Edge magnetotransport in graphene: A combined analytical and numerical study. *Annalen der Physik*, 527:723–736, 2015.
- [20] P. W. Anderson. Absence of diffusion in certain random lattices. *Phys. Rev.*, 109:1492, 1958.
- [21] F. Evers and A. D. Mirlin. Anderson transitions. *Rev. Mod. Phys.*, 80:1355, 2008.
- [22] A. Legendijk, B. van Tiggelen, and D. S. Wiersma. Fifty years of anderson localization. *Phys. Today*, 62:24, 2009.
- [23] T. Stegmann, O. Ujsághy, and D. E. Wolf. Localization under the effect of randomly distributed decoherence. *Eur. Phys. J. B*, 87:30, 2014.
- [24] L. V. Keldysh. Diagram technique for nonequilibrium processes. *Sov. Phys. JETP*, 20:1018, 1965.
- [25] L. P. Kadanoff and G. Baym. *Quantum Statistical Mechanics*. W. A. Benjamin, 1962.
- [26] S. Datta. *Quantum Transport: Atom to Transistor*. Cambridge University Press, 2005.
- [27] S. Datta. *Lessons from Nanoelectronics: A New Perspective on Transport*. World Scientific, 2018.
- [28] M. Di Ventra. *Electrical Transport in Nanoscale Systems*. Cambridge University Press, 2008.
- [29] C. Caroli, R. Combescot, P. Nozieres, and D. Nozieres. Direct calculation of the tunneling currents. *J. Phys. C*, 4:916, 1971.

- [30] A. Cresti, R. Farchioni, G. Grosso, and G. P. Parravicini. Keldysh-green function formalism for current profiles in mesoscopic systems. *Phys. Rev. B*, 68:075306, 2003.
- [31] R. Landauer. Spatial variation of currents and fields due to localized scatterers in metallic conduction. *IBM J. Res. Dev.*, 1:223, 1957.
- [32] R. Landauer. Electrical transport in open and closed systems. *Z. Phys. B*, 68:217, 1987.
- [33] E. N. Economou. *Green's Functions in Quantum Physics*. Springer, 2006.
- [34] T. Heikkilä. *The Physics of Nanoelectronics: Transport and Fluctuation Phenomena at Low Temperatures*. Oxford University Press, 2013.
- [35] T. Ihn. *Semiconductor Nanostructures*. Oxford University Press, 2010.
- [36] D. K. Ferry, S. M. Goodnick, and J. Bird. *Transport in Nanostructures*. Cambridge University Press, 2009.
- [37] J. C. Cuevas and E. Scheer. *Molecular Electronics: An Introduction to Theory and Experiment*. World Scientific, 2017.

Cosmología observacional con Redes Neuronales Artificiales

J. Alberto Vázquez,^{1, a} Ricardo Medel Esquivel,² and Isidro Gómez-Vargas²

¹*Instituto de Ciencias Físicas UNAM, Av. Universidad s/n,
Col. Chamilpa, Cuernavaca, Morelos, 62210, México.*

²*CICATA-Legaria, Instituto Politécnico Nacional, Ciudad de México, CP 11500, México.*

En estas memorias se describe el empleo de la inferencia Bayesiana en la cosmología observacional y cómo esta técnica se puede combinar con Redes Neuronales Artificiales para optimizar el tiempo de cómputo. En el transcurso del escrito, se pretende dar al lector un panorama muy general de la cosmología observacional, el torrente de datos que genera y la necesidad de incorporar algoritmos avanzados, tanto estadísticos como de inteligencia artificial, en este campo de investigación.

Palabras clave: cosmología observacional, estimación de parámetros, redes neuronales artificiales

I. INTRODUCCIÓN

Durante la última década, el rápido avance en el desarrollo de instrumentos observacionales altamente sofisticados, así como el incremento en el poder de cómputo han guiado al surgimiento de la *Cosmología de alta precisión*. En particular, experimentos desarrollados para medir las anisotropías observadas en la radiación cósmica de fondo, las distancias luminosas a partir de explosiones de Supernovas tipo Ia y la formación de estructura a gran escala. Por ejemplo, la Figura 1 muestra la distribución de varios millones de galaxias en las escalas más grandes del cosmos, cada punto representa una galaxia. Son estas escalas del Universo las que estudia la cosmología y, en particular, la cosmología observacional se enfoca en hacer uso de este tipo de información para validar los modelos teóricos que intentan explicar matemáticamente la formación, estructura y evolución del Universo.

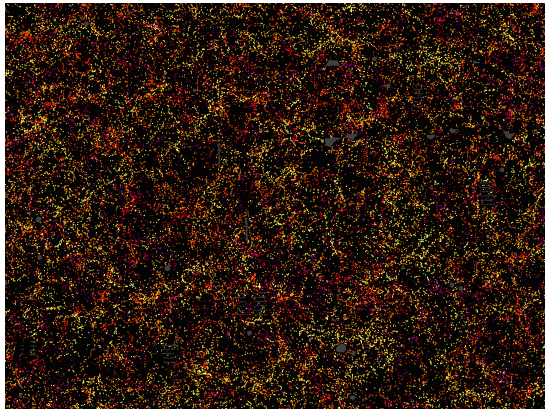


Figura 1. Distribución de varios millones de galaxias en las escalas más grandes del cosmos. Fuente: SDSS [1]

Esta imagen, la Figura 1, conjunta el trabajo proveniente del Sloan Digital Sky Survey (SDSS), un proyecto fundado en el año 2000, cuyo observatorio se encuentra

ubicado en Apache Point, Nuevo México, y que tiene como objetivo capturar imágenes y señales de las profundidades del Universo.

Así como este gran proyecto se enfoca en recabar información sobre la distribución de galaxias y oscilaciones acústicas de bariones desde la Tierra, existen también satélites, como COBE¹, WMAP² y Planck³, que proveen un laboratorio celeste para la cosmología y permiten estudiar otros fenómenos cósmicos desde el espacio.

Esta nueva etapa en el estudio del Universo permite combinar una serie de observaciones a diferentes escalas y, por tanto, probar nuestras teorías con mayor exactitud. Por si fuera poco, dentro de un par de años, el telescopio LSST⁴ fotografiará toda la bóveda celeste con la cámara más grande jamás construida (3.2 Gigapíxeles), y se estima que recopilará alrededor de 15 Terabytes por noche. Se vislumbra que la gestión y minería de datos del telescopio sean la parte técnica más difícil del proyecto. Por otro lado, el próximo radiotelescopio más grande del mundo, SKA⁵, coleccionará aproximadamente 14 exabytes de información por día, suficientes datos como para llenar 15 millones de iPods de 64 GB a diario; además, su computadora central tendrá el poder de, más o menos, 100 millones de unidades de procesamiento. La cosmología, así como otras áreas de la ciencia, ha entrado a una nueva etapa donde la cantidad de datos juega un papel imprescindible para el entendimiento del Universo.

Jim Gray, premio Turing en 1998 y desarrollador de la base de datos del SDSS, pronosticó el advenimiento de un cuarto paradigma en la ciencia. J. Gray mencionaba que el primer paradigma ocurrió hace mil años, cuando la ciencia se limitaba a ser descriptiva y empírica; el segundo tuvo su inicio con las investigaciones de Isaac Newton, al adquirir un carácter teórico; el tercero surgió en el siglo XX, con la incorporación de las simulaciones computacionales. El cuarto paradigma sería la ciencia de

¹ science.nasa.gov/missions/cobe

² map.gsfc.nasa.gov

³ <https://bit.ly/2lQJbfK>

⁴ www.lsst.org

⁵ <https://bit.ly/2kexoY5>

^a javazquez@icf.unam.mx

datos o *e-ciencia*, una combinación de ciencia teórica con análisis computacional de una cantidad exorbitante de datos [2].

La cosmología contemporánea, con su gran caudal de datos y la necesidad de analizarlos de manera exhaustiva, cristaliza en una realidad la visión de Jim Gray del cuarto paradigma científico.

En general, el análisis de los datos cosmológicos requiere de gran variedad de técnicas estadísticas y computacionales. Por ello, las redes neuronales artificiales y otros algoritmos de inteligencia artificial, en años recientes, son útiles dentro de la cosmología. En particular, en estas memorias, mencionaremos cómo se pueden usar las redes neuronales artificiales dentro de la estimación de parámetros cosmológicos.

II. LA INFERENCIA BAYESIANA

El modelo Λ CDM es considerado el “modelo estándar” en cosmología, puesto que es la parametrización más sencilla que toma en cuenta todos los aspectos y componentes fundamentales del Universo. Este modelo incluye una componente de propiedades exóticas para dar explicación a la actual expansión acelerada del Universo, comúnmente denominada energía oscura ($\Omega_{DE} = 1 - \Omega_m$), y cuya propuesta más simple es descrita a través de una constante cosmológica (Λ). Recientes observaciones indican que la energía oscura contribuye al 70% del contenido total del universo (Ω_{DE}), mientras que el 30% restante corresponde a la suma de la materia oscura fría (CDM) y materia ordinaria o bariónica Ω_m . La inferencia Bayesiana, en conjunto con las observaciones cosmológicas, nos permiten deducir los componentes del Universo [3]. Por ejemplo, la Figura 2 muestra dos tipos de gráficas que son muy importantes en el análisis de datos cosmológicos. La primera es conocida como distribución *posterior* 1-dimensional, donde se visualiza la región más probable de un parámetro, en este caso, de la densidad total de materia Ω_m . Por otro lado, la imagen de la derecha representa la distribución *posterior* 2-Dimensional, con las regiones de probabilidad de dos parámetros, en este caso, la densidad de materia Ω_m junto con la densidad bariónica Ω_b ; la parte interna de estos elipses representa la región de confianza 1σ (68.3%), mientras que la externa corresponde a 2σ (95.4%). Este tipo de análisis se obtiene después de realizar un proceso de inferencia bayesiana, que se basa en datos observacionales así como en un modelo teórico dado.

El teorema de Bayes es la base de la inferencia bayesiana y, en la prueba de modelos teóricos, tiene el siguiente aspecto:

$$P(\theta|D, H) = \frac{P(D|\theta, H)P(\theta|H)}{P(D|H)},$$

donde D representa el conjunto de datos observacionales, H es la hipótesis (el modelo) y θ el conjunto de

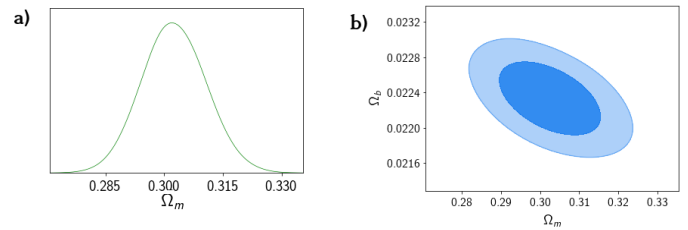


Figura 2. Distribución *posterior* sobre el contenido del Universo: a)1-Dimensional, b)2-Dimensional.

parámetros (un modelo que describe una línea, por ejemplo, tendría como parámetros libres la pendiente de una recta y su ordenada al origen). Para el caso en cuestión, un ejemplo de análisis en cosmología, estas cantidades describen los siguientes aspectos:

- D - Datos/Observaciones cosmológicas, i.e. aquellos provenientes del SDSS.
- H - Modelo cosmológico, i.e. modelo Λ CDM.
- θ - Conjunto de parámetros que mejor describen las observaciones, i.e. contenido de energía oscura Ω_{DE} .

Además, la probabilidad *previa* $P(\theta|H)$ ⁶ representa nuestro conocimiento de los parámetros θ antes de considerar los datos observables. Esta probabilidad se modifica a través de la *verosimilitud* $P(D|\theta, H)$ al incluir datos experimentales D . El objetivo final es obtener la probabilidad *posterior* $P(\theta|D, H)$, que representa el estado de nuestro conocimiento de los parámetros del modelo al considerar la información de los datos. La constante de normalización $P(D|H)$, o *evidencia bayesiana*, es el promedio de la *verosimilitud* sobre la probabilidad *previa*:

$$P(D|H) = \int d^N \theta P(D|\theta, H)P(\theta|H),$$

donde N es la dimensión del espacio de parámetros. Esta cantidad, fundamental en la comparación de modelos, es una constante que puede omitirse en la estimación de parámetros.

La *verosimilitud* se determina al suponer una distribución estadística para los datos D , la probabilidad *previa* se determina a partir de información conocida con antelación. Cuando se calcula la distribución *posterior*, se conocen los valores más probables de los parámetros basados en el modelo teórico y en los datos observacionales.

En el cálculo de las integrales de alto grado de complejidad, que surgen de manera natural en la inferencia bayesiana, se utilizan métodos numéricos especiales, conocidos como Métodos de Monte Carlo vía Cadenas de Markov (MCMC, por sus siglas en inglés).

⁶ Las probabilidades involucradas son condicionales.

Explicado de manera simplificada, el método consiste en realizar un muestreo de las distribuciones de probabilidad bajo la integral, es decir, en tomar aleatoriamente una gran cantidad M de valores posibles en el espacio N -dimensional de parámetros acordes a la densidad de probabilidad expresada por la *verosimilitud*, de tal modo que se satisfagan las condiciones del Teorema Ergódico, que es una generalización de la ley fuerte de los grandes números válida para cadenas de Markov. Este muestreo, entonces, permite aproximar la integral como una sumatoria:

$$\int d^N \theta P(D|\theta, H) P(\theta|H) \approx \frac{1}{M} \sum_{i=1}^M \mathcal{L}(\theta_i),$$

donde $\mathcal{L}(\theta)$ es una función construida a partir de la *verosimilitud* y depende de la cantidad que se desea inferir.

Las mayores dificultades técnicas de los métodos MCMC radican en el proceso de muestreo, pues es deseable explorar todo el espacio de búsqueda, es decir, ser capaces de elegir al azar cualquier valor posible en el espacio de parámetros. Sin embargo, las cadenas de Markov son secuencias de valores aleatorios con la característica distintiva de que cada valor depende exclusivamente del valor anterior; característica que finalmente induce la convergencia de la cadena de Markov a una distribución de probabilidad específica pero también reduce la región de exploración aleatoria.

Tales dificultades técnicas han estimulado la mejora de los métodos MCMC y la búsqueda de nuevos esquemas de exploración estocástica, como el muestreo anidado, siendo en la actualidad un campo de investigación muy activo.

En el caso particular de la cosmología, dado que los modelos cosmológicos suelen tener al menos seis parámetros, diversos grupos de investigación han desarrollado códigos especializados que agrupan varios algoritmos para hacer la estimación de parámetros, por ejemplo: CosmoMC [4], CosmoSIS [5] y SimpleMC [6].

Al lector que desee ahondar en la inferencia bayesiana se le recomienda la Ref. [7] como texto de divulgación introductorio, la Ref. [8] para una introducción con mayores detalles técnicos, la Ref. [9] si se desea un tratamiento más formal, [3] es un buen tratamiento en el contexto cosmológico y en [10] se puede encontrar un texto de divulgación, similar al presente, con mayor énfasis en determinado tipo de algoritmos que permiten la inferencia bayesiana y comparación de modelos.

III. REDES NEURONALES ARTIFICIALES

En la sección anterior, se mencionó que los algoritmos involucrados en la inferencia bayesiana de un modelo cosmológico tienen que explorar un espacio de parámetros con muchas dimensiones. Las integrales involucradas para el cálculo de la función de densidad de probabilidad de *verosimilitud* y de la *evidencia bayesiana* se tornan muy

complejas. Una alternativa para optimizar estos cálculos es mediante una red neuronal artificial.

La arquitectura más simple de las redes neuronales artificiales, pertenece a la clase de algoritmos del *aprendizaje automático* denominados supervisados.

En la Figura 3 se esquematiza, a grandes rasgos, las partes que componen a una red neuronal artificial simple. Cada neurona recibe datos de entrada de todas las neuronas que le anteceden. La suma ponderada es un valor que contiene la información de los valores de entrada multiplicada por los pesos w que se les asignan. Después, esta suma es el parámetro de una función no lineal de activación. La salida final de todas las neuronas interconectadas, es la que se compara con el valor esperado y con la cual se calcula el error que permite entrenar a la red neuronal.

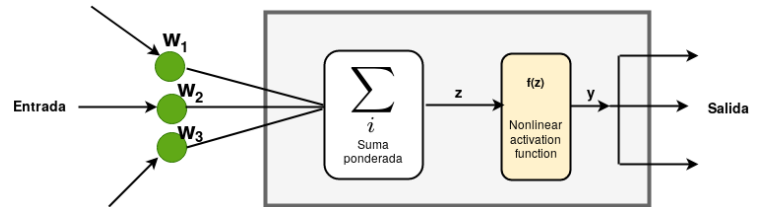


Figura 3. Esquema simplificado de una red neuronal artificial

El entrenamiento consiste en minimizar una función objetivo, llamada función de costo, que mide la diferencia entre las predicciones (aleatorias al inicio del proceso) y los valores reales conocidos de antemano. Esta minimización del error (medido por la función de costo) se efectúa después de varias iteraciones donde se aplican los algoritmos del *descenso del gradiente* y *retropropagación* (ver Figura 4).

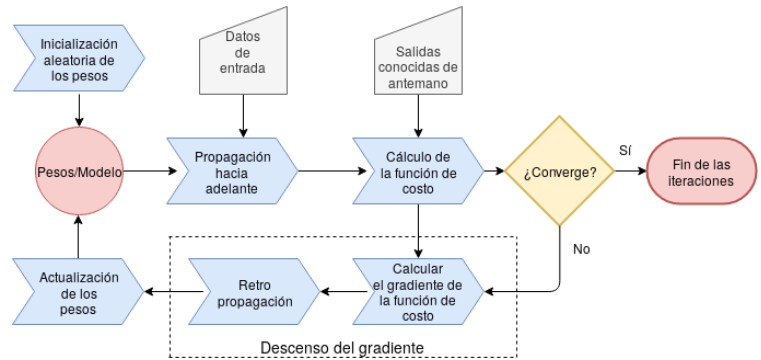


Figura 4. Entrenamiento de una red neuronal

En la Fig. 5 se puede observar un perceptrón multicapa, donde las neuronas se ordenan en paredes, y todas las neuronas de una capa reciben información de las que les preceden.

Debido a que la inferencia bayesiana es un método muy bien estructurado, las redes neuronales se pueden entrenar mediante el muestreo del espacio de parámetros y con los cálculos de las integrales que el algoritmo en

función esté realizando. Cabe mencionar, que el ajuste de parámetros de un modelo cosmológico puede requerir más de 10^6 iteraciones, donde se puede invertir el 80 % de las iteraciones en entrenar una red neuronal como la de la Fig. 5 y en las restantes utilizar la red neuronal entrenada para hacer el cálculo de las complejas integrales, y por tanto disminuir considerablemente el tiempo de computo.

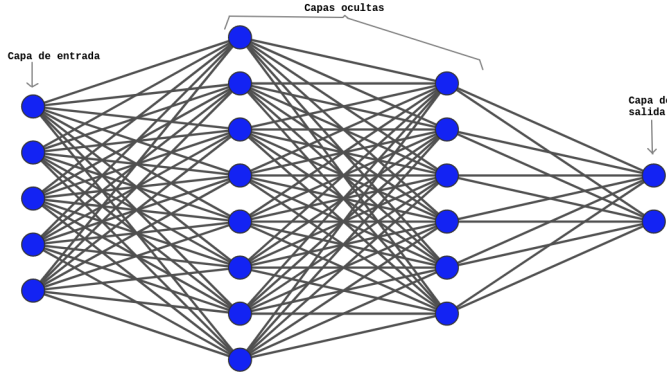


Figura 5. Red Neuronal Artificial con dos capas ocultas

IV. EJEMPLO

A modo de ejemplo, estimamos los parámetros del modelo cosmológico LCDM con el algoritmo conocido como muestreo anidado elipsoidal [11] y, también, con ese mismo algoritmo complementado mediante redes neuronales artificiales. La implementación que hemos realizado fue incorporada dentro del código SimpleMC.

Los datos observacionales utilizados fueron SN-IA, cronómetros cósmicos, oscilaciones acústicas de bariones y datos de la radiación cósmica de fondo obtenidos por el satélite Planck. En la Figura 6, se pueden apreciar los resultados obtenidos mediante ambos métodos.

Además, en la Tabla I, se comparan los tiempos que consumieron ambos métodos, así como sus cálculos de la evidencia bayesiana. Esta diferencia de tiempo, a pesar de parecer pequeña ($\sim 13\%$), corresponderá a una mejora considerable cuando se consideren modelos cosmológicos más y una cantidad elevada de datos, como la que esperamos en los siguientes años.

	sin RNA	con RNA
tiempo (minutos)	3.194	2.819
$\log(Evidencia)$	$-35,611 \pm 0,233$	$-35,67 \pm 0,231$

Tabla I. Tiempos de ejecución y evidencias bayesianas

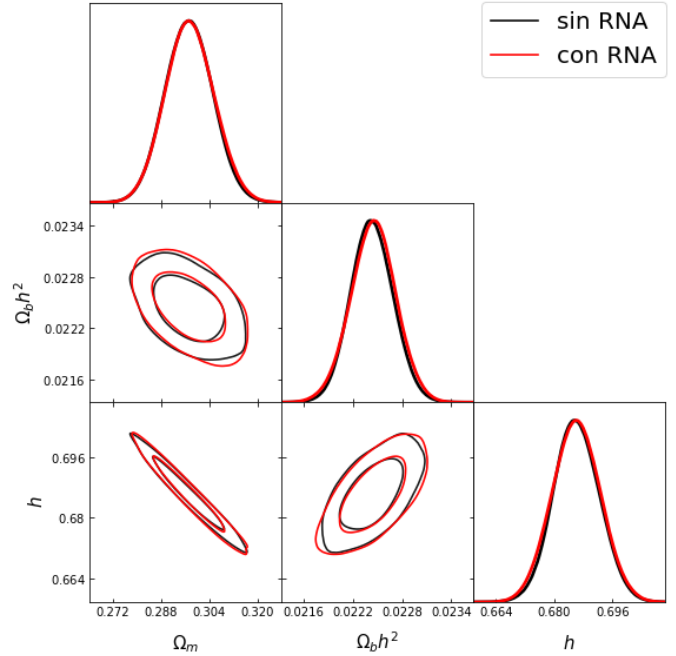


Figura 6. Comparación de resultados

V. CONCLUSIONES

La incorporación de redes neuronales dentro de la estimación de parámetros cosmológicos, permite disminuir el tiempo computacional al aprender a calcular las integrales de las funciones de *verosimilitud*. Además, se conserva la precisión tanto en la estimación de los valores más probables de los parámetros como en el cálculo de la evidencia bayesiana.

AGRADECIMIENTOS

J.A.V. agradece el apoyo proporcionado a los proyectos FOSEC SEP-CONACYT Investigación Básica A1-S-21925 y UNAM-PAPIIT IA102219. I. G. V y R. M. E. agradecen el apoyo brindado por el CONACYT a través de sendas Becas de Posgrado.

[1] Astronomers map a record-breaking 1.2 million galaxies to study the properties of dark energy. <https://bit.ly/21MDQ92>. Accessed: 2019-09-12.

[2] T Hey, S Tansley, and K Tolle. Jim gray sobre la eficiencia: un método científico transformado. In *El cuarto paradigma. Descubrimiento científico intensivo en datos*, chapter 1, pages 17–23. UAM, México, 2014.

- [3] Luis E Padilla, Luis O Tellez, Luis A Escamilla, and J Alberto Vazquez. Cosmological parameter inference with bayesian statistics. *arXiv preprint arXiv:1903.11127*, 2019.
- [4] Antony Lewis. Efficient sampling of fast and slow cosmological parameters. *Physical Review D*, 87(10):103529, 2013.
- [5] Joe Zuntz, Marc Paterno, Elise Jennings, Douglas Rudd, Alessandro Manzotti, Scott Dodelson, Sarah Bridle, Saba Sehrish, and James Kowalkowski. Cosmosis: modular cosmological parameter estimation. *Astronomy and Computing*, 12:45–59, 2015.
- [6] Éric Aubourg, Stephen Bailey, Julian E Bautista, Florian Beutler, Vaishali Bhardwaj, Dmitry Bizyaev, Michael Blanton, Michael Blomqvist, Adam S Bolton, Jo Bovy, et al. Cosmological implications of baryon acoustic oscillation measurements. *Physical Review D*, 92(12):123516, 2015.
- [7] P Castro Ortega. *El Teorema de Bayes. Aproximándonos a la verdad*. RBA coleccionables, España, 2017.
- [8] R Medel Esquivel, I Gómez-Vargas, JA Vázquez, and R. García Salcedo. An introduction to markov chain monte carlo. *Boletín de Estadística e Investigación Operativa*, 2019. Accepted but not published yet.
- [9] Devinderjit Sivia and John Skilling. *Data analysis: a Bayesian tutorial*. OUP Oxford, 2006.
- [10] I Gómez-Vargas, R. Medel Esquivel, R. García Salcedo, and JA Vázquez. Una aplicación de las redes neuronales en la cosmología. *Komputer Sapiens*, 11-2, 2019. DOI: 10.13140/RG.2.2.31865.52321. Accepted but not published yet.
- [11] Pia Mukherjee, David Parkinson, and Andrew R Liddle. A nested sampling algorithm for cosmological model selection. *The Astrophysical Journal Letters*, 638(2):L51, 2006.

Redes Reguladoras Complejas en la Medicina y la Biología.

María Barrera
*Posgrado en Ciencias Biomédicas
Instituto de Ecología
Universidad Nacional Autónoma de México*

Carlos Villarreal
*Instituto de Física
Universidad Nacional Autónoma de México*

I. INTRODUCCIÓN

El estilo de vida moderno caracterizado por el consumo excesivo de carbohidratos, el sedentarismo, la obesidad, la acción de contaminantes o el estrés emocional, entre otros factores. Se sabe ahora que dicho estilo de vida está asociado al desarrollo de enfermedades con una sintomatología muy diferente, tales como la diabetes tipo 2, los accidentes cardiovasculares, los derrames cerebrales o el cáncer de colon. Si bien estos padecimientos poseen un carácter multifactorial, involucran un denominador común: una respuesta inmune expresada como un estado inflamatorio crónico de bajo nivel [1].

Normalmente, ante la invasión de un agente patógeno las células del sistema inmune pueden desarrollar diferentes clases de respuestas efectoras, que pueden ser del tipo celular, humoral o inflamatorio. En particular, la respuesta inflamatoria involucra la eliminación rápida de virus, bacterias y hongos que expresan motivos moleculares característicos: azúcares de membrana, proteínas flagelares, fragmentos de ARN, etc., los cuales son reconocidos por las diversos tipos de células inmunitarias encargadas de esta respuesta; sin embargo, el exceso de azúcares y grasas circulantes en el organismo conduce a la alteración de los procesos metabólicos involucrados en la transformación de nutrientes en energía, así como a la liberación de agentes químicos pro-inflamatorios que se traduce en un deterioro paulatino de diversos tejidos y a la generación a largo plazo de enfermedades crónico degenerativas.

La naturaleza multifactorial de este tipo de enfermedades sugiere que una herramienta ideal para su modelado son las redes complejas, las cuales permitan describir las interacciones mutuas de las diversas componentes celulares, moleculares y ambientales involucrados en su desarrollo. La observación de que una respuesta inmune crónica de bajo nivel induce un deterioro gradual de la funcionalidad celular y enfermedades concomitantes, conducen a la propuesta de que éstas pueden entenderse en términos de transiciones de fase entre un estado relativamente sano, caracterizado por una homeostasis celular alterada, pero capaz de compensar su progresivo deterioro funcional, hacia otros estadios asociados con el agotamiento de dicha capacidad compensatoria y la expresión manifiesta de estos padecimientos. En el lenguaje de los sistemas complejos, las diferentes fases tienen su origen en un conflicto entre fuerzas efectivas que favorecen la expresión de un fenómeno con otras que se contraponen al mismo, y cuya resolución conduce a una síntesis manifiesta por la emergencia de patrones de expresión estacionarios. Dicha conjetura puede

formalizarse, como se verá posteriormente, recurriendo a la teoría de sistemas cooperativos, transiciones de fase y autoorganización [2].

En lo que sigue, se presentará una exposición de herramientas conceptuales y matemáticas involucradas en la teoría de sistemas complejos en la descripción de las enfermedades con origen inflamatorio. Primero se abordará el estudio de la dinámica inmunitaria en términos de procesos de autoorganización y transiciones de fase fuera del equilibrio. A continuación, se discutirá el abordaje de estas enfermedades mediante redes reguladoras conformadas por componentes celulares, moleculares y ambientales involucrados en su desarrollo. Para ello se introducirán conceptos de las redes con reglas discretas booleanas. Posteriormente, se generalizará esta descripción a redes complejas con interacciones continuas descritas por la lógica difusa.

II. TRANSICIONES DE FASE LEJOS DEL EQUILIBRIO Y AUTOORGANIZACIÓN

La formación espontánea de estructuras organizadas puede explicarse mediante la teoría de fenómenos cooperativos, auto-organización y transiciones de fase fuera del equilibrio [2]. Los fenómenos cooperativos surgen de las interacciones no lineales de un gran número de subsistemas elementales que conducen a la aparición de patrones o fases organizadas. Esta teoría ha sido empleada en el análisis de múltiples sistemas que desarrollan autoorganización: la magnetización en ferromagnetos, la generación de luz láser, los fenómenos de superfluidez y superconductividad, o la manifestación del campo de Higgs de las interacciones fundamentales. El concepto 'fuera del equilibrio' se refiere a que estos patrones revisten un equilibrio estacionario, pero no necesariamente de equilibrio termodinámico. La teoría se basa en dos conceptos principales, la existencia de parámetros de orden y de control. Los parámetros de orden son aquellas variables que se identifican con la emergencia de diferentes fases de organización, entre las cuales puede transitar un sistema dependiendo de los valores específicos de los parámetros de control. La identificación de los parámetros de orden en una teoría no es directa, y requiere un análisis cuidadoso de la relevancia de las diferentes variables de un sistema, por ejemplo, aquellas asociadas a elementos con un alto grado de conectividad con otros elementos. Por otro lado, los parámetros de control están asociados al entorno, es decir, a las restricciones a las que está sujeto

el sistema.

En la descripción de la evolución de una enfermedad con origen inflamatorio crónico debe tomarse en cuenta, como ya se ha mencionado anteriormente, que ante la presencia de agentes patógenos o el exceso de nutrientes tales como grasas o azúcares, las células del sistema inmune pueden desarrollar una respuesta efectora del tipo inflamatorio. En forma concomitante, ocurre una respuesta reguladora que tiende a inhibir efectos perjudiciales derivados de una acción efectora exacerbada. La regulación mutua de los diferentes tipos de respuestas requiere de una comunicación transmitida mediante la secreción de mensajeros químicos (quimiocinas) cuyo efecto depende, tanto de la concentración de las células excretoras y receptoras, como de su afinidad hacia receptores celulares específicos. Denotemos por N^+ al número de células desarrolla una respuesta efectora del tipo inflamatorio, y por N^- al número de células que genera una respuesta reguladora, de modo que $N^+ + N^- = N$. Si introducimos ahora las fracciones relativas, $n^+ = N^+/N$ y $n^- = N^-/N$, con $n^+ + n^- = 1$, las consideraciones anteriores nos permiten suponer que la comunicación celular induce acciones colectivas que favorecen o inhiben la respuesta inmune efectora, las cuales pueden representarse por términos proporcionales a $(n^+)^2 n^-$, y $-(n^-)^2 n^+$, respectivamente. Asimismo, supondremos que en ausencia de estas interacciones decae a una tasa a . En tal caso, la dinámica de la fracción efectora puede describirse por la ecuación

$$\frac{dn^+}{dt} = c(n^+)^2 n^- - b(n^-)^2 n^+ - an^+ \quad (1)$$

$$= -\gamma(n^+)^3 + \beta(n^+)^2 - \alpha n^+ \equiv F(n^+), \quad (2)$$

en donde la Ec.(2) se obtiene sustituyendo la condición $n^- = 1 - n^+$ en la Ec.(1). En lo que sigue, es conveniente redefinir la tasa de decaimiento efectiva $\alpha \rightarrow \alpha - \alpha_c$, e introduciré la notación $n \equiv n^+$. En tal caso, los estados de equilibrio estacionario del sistema satisfacen la condición $dn/dt = 0$, con soluciones $n = 0$ y

$$n = \frac{\beta}{2\gamma} \pm \frac{1}{2\gamma} [\beta^2 - 4\gamma(\alpha - \alpha_c)]^{1/2}. \quad (3)$$

En general, todas la raíces son admisibles; sin embargo, en el presente problema se debe cumplir que $n \geq 0$, por lo que se deben considerar sólo las soluciones con la raíz positiva; el resto de los parámetros debe satisfacer las relaciones: $\beta \geq 0$, $\gamma > 0$, y $\alpha \leq \beta^2/4\gamma$.

En forma análoga a los sistemas mecánicos, en esta descripción se puede introducir un potencial de configuraciones en el espacio de estados, $V(n)$, tal que $F(n) = -\partial V(n)/\partial n$. Podemos visualizar este potencial como un paisaje con colinas y valles; es fácil ver que la condición de equilibrio estacionario implica que dichos estados residen en los valles del mismo. Las interacciones involucradas en la Ec. (1) están descritas por el potencial:

$$V(n) = \frac{1}{2}(\alpha - \alpha_c)n^2 - \frac{1}{3}\beta n^3 + \frac{1}{4}\gamma n^4. \quad (4)$$

Examinemos primero el caso en que $\beta = 0$. Entonces $V(n)$ se transforma en un potencial 'cuártico' cuyos valores mínimos

están determinados por el parámetro α : si $\alpha > \alpha_c$, el mínimo se encuentra en $n = 0$, mientras que para $\alpha < \alpha_c$, el mínimo se encuentra en $n = (|\alpha - \alpha_c|/\gamma)^{1/2} \equiv n_e$. Vemos que α actúa como un parámetro de control que determina una transición de fase de un estado inactivo a otro activo definido por el valor finito del parámetro de orden n . En el modelo de infección por el VIH, este proceso se identifica con una transición de una respuesta inmune de tolerancia al virus, a una respuesta efectora exacerbada que da lugar a una aceleración de la reproducción viral y sus efectos perjudiciales. Estudios experimentales confirman que la activación crónica de la respuesta inmune conduce a un deterioro gradual del tejido linfático, que es donde se aloja el 99% del virus, y que dicho deterioro correlaciona con el advenimiento del SIDA [3].

Una fenomenología similar ocurre si consideramos que ahora β tiene un valor finito; sin embargo, el carácter de la transición se modifica, ya que en el primer caso la transición entre estados ocurre en forma continua, mientras que en el segundo ocurre en forma discontinua. De acuerdo con una clasificación estándar de transiciones de fase, en el primer caso ésta es de segundo orden, con $V'(0) = V'(n_e)$, pero $V''(0) \neq V''(n_e)$, mientras que en el primero es de primer orden, con $V'(0) \neq V'(n_e)$.

III. REDES REGULADORAS BOOLEANAS Y PAISAJE EPIGENÉTICO

Una concepción semejante a la desarrollada por Germinal para el estudio de la complejidad biológica fue elaborada en forma independiente por C. H. Waddington en 1957, con la introducción del concepto metafórico del paisaje epigenético [4]. Esta idea tiene como fin entender los procesos reguladores de genes involucrados en el desarrollo celular, el cual puede visualizarse como una bola rodando en un paisaje formado por picos y valles. Siguiendo su trayectoria, la pelota finalmente puede caer en un valle, representando su posición final un destino o fenotipo celular. Este concepto fue formalizado posteriormente por S. A. Kauffman y L. Glass [5, 6], quienes estudiaron el comportamiento de grandes redes de 'genes' binarios interconectados, estableciendo los principios del modelado booleano de procesos de diferenciación celular.

El modelado matemático basado en redes con reglas booleanas [7–9] proporciona información cualitativa significativa sobre la topología básica de las relaciones que determinan alternativas destinos celulares y pueden usarse para el análisis de circuitos biológicos sin requerir valores explícitos de los parámetros involucrados en la red. En este tipo de enfoque, los nodos de la red representan genes, factores de transcripción, proteínas moduladoras de vías de señalización, factores ambientales, etc., y enlaces que representan regulaciones positivas o negativas entre pares de nodos. La variable de estado de cada nodo toma un valor discreto de 0 (inhibido o inactivo) o 1 (expresado o activo). El estado de cada nodo en el momento $t + 1$ se especifica mediante un mapeo dinámico que depende del estado de sus reguladores en un momento anterior t :

$$q_k(t + 1) = F_k(q_1(t), \dots, q_n(t)), \quad (5)$$

donde F_k es una función discreta que representa una proposición lógica, definida por una regla booleana, es decir, satisface la axiomática establecida por G. Boole, quien desarrolló un sistema de reglas que permitían expresar, manipular y simplificar problemas lógicos cuyos argumentos admiten dos estados (verdadero o falso). Si denotamos los conectivos lógicos como $\text{and} = \wedge$, $\text{or} = \vee$, $\text{not} = \neg$, los axiomas de Boole para la conjunción lógica son:

$$\begin{array}{ll} a \wedge (b \wedge c) = (a \wedge b) \wedge c & \text{asociatividad} \\ a \wedge b = b \wedge a & \text{conmutatividad} \\ a \wedge (a \vee b) = a & \text{absorción} \\ a \wedge (b \vee c) = (a \wedge b) \vee (a \wedge c) & \text{distributividad} \\ a \wedge \neg a = \emptyset & \text{complemento} \end{array}$$

Los axiomas correspondientes a la disyunción lógica pueden derivarse directamente de las leyes de De Morgan, que definen la negación de las proposiciones booleanas compuestas

$$\neg(a \wedge b) = (\neg a) \vee (\neg b),$$

$$\neg(a \vee b) = (\neg a) \wedge (\neg b).$$

La dinámica inducida por el mapeo booleano se determina en forma completa una vez que se especifica un conjunto de valores de expresión inicial de los componentes de la red. A partir de un conjunto inicial dado, los nodos actualizan iterativamente su valor en sincrónicamente, hasta que finalmente alcanzan un estado estable determinado por la condición $q_k(t+1) = q_k(t)$. Esta última condición especifica un atractor de punto fijo. Alternativamente, en la simulación de algunas redes reguladoras puede, resultar atractores cíclicos asociados a la condición $q_k(t+N) = q_k(t)$, en donde el número entero N denota el período del atractor. Esta clase de atractores describe comportamientos oscilatorios y son determinados por al menos un circuito de retroalimentación negativa en la topología de la red, equivalente a un número impar de interacciones del tipo inhibitorio. En el análisis de la dinámica, un papel especial es jugado por circuitos de pares de nodos mutuamente inhibitorios $(q_i \neg q_k) \wedge (q_k \neg q_i)$, lo cuales funcionan como 'switches' que determinan vías alternativas de señalización de la red, es decir, conforman núcleos de bifurcación que conducen a diferentes atractores.

En contraste con el esquema sincrónico, también se puede considerar un enfoque iterativo asíncrono, en el que cada nodo actualiza su valor de acuerdo a: $q_k(t + \tau_k) = F_k(q_1(t), \dots, q_n(t))$, donde τ_k es un tiempo de actualización característico. En este caso, la condición de punto fijo, $q_k(t + \tau_k) = q_k(t)$, conduce exactamente al mismo conjunto de relaciones obtenidas en el método síncrono, y por lo tanto, los conjuntos de atractores de punto fijo coinciden en ambos esquemas. Sin embargo, la condición para que exista un atractor periódico, $q_k(t + \tau_k) = q_k(t + N)$, implica que los tiempos de actualización de los distintos nodos deben satisfacer la relación $\tau_i/\tau_k = N$, para todo i, k , de tal manera que el número de tractors cíclicos se reduce respecto del caso síncrono. Asimismo, el tamaño de las cuencas de atracción puede diferir también. De hecho, en el marco asíncrono, el tamaño

de las cuencas depende de las diferentes opciones de los pasos de iteración τ_k .

La dinámica general de un modelo puede evaluarse rastreando las trayectorias de todas las configuraciones iniciales posibles en el espacio de estados hacia los atractores; el tamaño de este espacio está dado por $\Omega = 2^n$ donde n es el número de nodos en la red. En general, cada punto fijo o atractor cíclico se puede alcanzar desde un número ω de diferentes condiciones iniciales. El parámetro ω denota el tamaño de la cuenca de atracción que puede visualizarse como el área de una región del paisaje epigenético. En consecuencia, la probabilidad de que se exprese un estado estacionario viene dada por una proporción de áreas, $p = \omega/\Omega$. El comportamiento general del sistema y la cantidad de atractores de una red reguladora booleana depende tanto de sus características topológicas, como de la cantidad de sus componentes y el grado de interconectividad entre ellas. Se reconoce ahora que las redes biológicas son sistemas libres de escala, lo que significa que sus nodos muestran una gran diversidad en el número de aristas, incluyendo pocos elementos con muchos enlaces y muchos elementos con pocos enlaces. La invariancia ante transformaciones de escala proporciona, entre otros atributos: robustez de la red, mayor eficiencia en la propagación de la información y la propiedad de que el número de atractores es casi independiente del número de nodos [10].

Los modelos booleanos han sido útiles para integrar datos publicados independientemente de diferentes circuitos moleculares involucrados en la especificación celular, para investigar cómo estos circuitos orquestan la diferenciación, y para generar nuevas hipótesis sobre interacciones no contempladas o estados celulares de diferenciación alternativos. Con base en estos modelos ha sido posible dilucidar las relaciones epigenéticas en plantas, donde el sistema paradigmático es la *Arabidopsis thaliana*. En este sistema se han realizado investigaciones sobre la generación de los órganos florales y la raíz [8, 11–13], la emergencia de patrones celulares y su regulación por fuerzas mecánicas y químicas [14], así como el comportamiento de los ciclos celulares de la planta [15]. En el área de la Biomedicina, se ha estudiado la influencia de un microambiente celular inflamatorio en el desarrollo de la leucemia [16], o el problema de la diferenciación y plasticidad celular de células del sistema inmune y su relación con las enfermedades crónico-degenerativas [17, 18]. Una revisión de las aplicaciones tanto de la teoría de redes reguladoras complejas para el estudio de procesos del desarrollo y plasticidad de células del sistema inmune en la generación de este tipo de enfermedades puede consultarse en la Ref. [9].

IV. REDES REGULADORAS Y LÓGICA DIFUSA

El estudio de las enfermedades crónicas ha influido fuertemente en la comprensión de cómo ligeros cambios derivan en una perturbación sistémica del comportamiento de sistemas biológicos complejos. Si se quisiera simular, por ejemplo, la forma en que la acumulación progresiva de factores proinflamatorios en el tracto intestinal perturban las proporciones de poblaciones de células T-CD4, la naturaleza discreta de

los valores de verdad involucrados los modelos booleanos implicaría un uso muy limitado para investigar las etapas transitorias entre el atractor sano y un atractor patológico [9]. Un enfoque más realista debe considerar que los niveles de expresión, las concentraciones y los parámetros de los sistemas biológicos pueden tomar cualquier valor dentro de un rango continuo limitado sólo por restricciones de funcionalidad [9, 19–21]. La traducción de las reglas booleanas interactivas al dominio continuo se puede lograr considerando un enfoque basado en la lógica difusa [22–25]. Esta teoría tiene como propósito el proporcionar una base formal para el razonamiento aproximado, y cuenta con aplicaciones tanto en la ingeniería de control, como en ciencias físicas, biomédicas y conductuales. Se caracteriza por un enfoque gradualista, de modo que el grado en que un objeto exhibe una propiedad determinada se especifica mediante una función de pertenencia (o característica) $\mu[w_k]$, con valores de verdad que van desde la falsedad total, $\mu[w_k] = 0$, hasta totalmente cierto, $\mu[w_k] = 1$. Por ejemplo, la propiedad de 'ser una persona talentosa' implica que hay un conjunto de personas que comparten ciertas características sin límite definido.

La lógica difusa satisface un axiomático similar al involucrado en la lógica booleana, excepto por el principio de identidad, lo que significa que el principio de no contradicción no es válido. Por lo tanto, aunque parezca paradójico, una proposición w y su negación $1 - w$ pueden ser simultáneamente verdaderas. En este contexto, la afirmación 'no era un tipo bondadoso, pero tampoco era malo' tiene un significado en la teoría del lenguaje. En los sistemas biológicos, las proposiciones difusas pueden describir casos en los que una célula muestra un patrón de expresión intermedio que no necesariamente pertenece a un fenotipo específico. La ausencia de no contradicción se expresa formalmente mediante la ecuación $w = 1 - w$, con la solución $w = 1/2$ [23]. Se deduce que el valor $w \equiv w_{thr} = 1/2$ puede interpretarse como un umbral entre la falsedad y la verdad (de lo anterior, se desprende que un discurso basado en medias verdades puede conducir a demostrar cualquier afirmación). Similarmente al enfoque booleano, en el régimen continuo las interacciones reguladoras de la red se caracterizan por proposiciones de lógica difusa denotadas aquí como $w_k[q_1(t), \dots, q_n(t)]$ y que se infieren de observaciones experimentales o son sugeridas por requisitos de consistencia interna. De hecho, un esquema de traducción del escenario discreto al continuo puede implementarse de manera directa mediante la sustitución de los conectivos booleanos and (\wedge), or (\vee) y not (\neg), por sus contrapartes difusas. Sin embargo, la definición de conectivos difusos no es única, y se han propuesto varias alternativas diferentes que no son completamente equivalentes. En la siguiente tabla presentamos la propuesta original de Zadeh [22] y un esquema de tipo 'probabilístico' [25]. Ambos esquemas satisfacen la axiomática booleana modificada discutida anteriormente. El segundo esquema muestra las mismas propiedades que las distribuciones de probabilidad conjuntas para variables independientes, de modo que las proposiciones difusas pueden traducirse directamente en enunciados de inferencia probabilística:

Boolean	Zadeh	'Probabilista'
q and p	$\min[q, p]$	$q \cdot p$
q or p	$\max[q, p]$	$q + p - q \cdot p$
not p	$1 - p$	$1 - p$

Un ejemplo de traducción del marco booleano al difuso está dado por la expresión:

$$w[p, q, r] = (q \vee p) \wedge \neg r \rightarrow (q + p - q \cdot p) \cdot (1 - r).$$

Las proposiciones lógicas continuas se pueden utilizar para construir una expresión explícita de la función característica $\mu[w_k]$. En el enfoque booleano discreto, esta función sería equivalente a una función Θ del tipo escalón:

$$\mu[w_k] \rightarrow \Theta[w_k - 1/2] = \begin{cases} 0 & \text{si } w_k < 1/2; \\ 1/2 & \text{si } w_k = 1/2; \\ 1 & \text{si } w_k > 1/2. \end{cases}$$

En el enfoque continuo, el grado de verdad que alcanza una proposición lógica puede representarse por una función característica con una estructura sigmoideal que cambia gradualmente de un valor nulo a un valor unitario. Muchas funciones comparten esta propiedad. Una expresión empleada en diferentes aplicaciones de lógica difusa en biología de sistemas es la función logística:

$$\mu[w_k] = \frac{1}{1 + \exp[-\beta(w_k(q_1, \dots, q_n) - w_{thr})]} \quad (6)$$

Aquí, w_{thr} es un valor umbral tal que si $w_k > w_{thr}$, w_k tiende a ser verdadero (o expresado). Por lo general, $w_{thr} = 1/2$. El parámetro β es una tasa de saturación que mide el ritmo del tránsito de un estado no expresado a uno expresado. Este ritmo es gradual para valores pequeños de β , y abrupto para grandes valores. En el caso $\beta \gg 1$, $\mu[w_k] \rightarrow \Theta[w_k - w_{thr}]$.

La dinámica discreta implicada por la Ec. (1) puede generalizarse al continuo mediante la introducción de un conjunto de ecuaciones diferenciales ordinarias para la evolución de los niveles de expresión $q_k(t)$ de los componentes de la red:

$$\frac{dq_k}{dt} = \mu[w_k(q_1, \dots, q_n)] - \alpha_k q_k, \quad (7)$$

donde $\mu[w_k]$ expresa una realización continua de la regla booleana w_k , mientras que α_k es una tasa de decaimiento. En este esquema, los estados de equilibrio del sistema están definidos por la condición de estado estable $dq_k/dt = 0$, que conduce a

$$q_k^s = \frac{1}{\alpha_k} \mu[w_k(q_1^s, \dots, q_n^s)], \quad (8)$$

donde el superíndice s denota el valor de estado estacionario. Una consecuencia directa de esto es que el nivel de expresión del nodo k depende en gran medida de su tasa de decaimiento. En el caso $\alpha_k > 1$, un nodo se sub-expresa con respecto al valor alcanzado por $\alpha_k = 1$; en particular, para $\alpha_k \gg 1$, la expresión de ese nodo se inhibirá por completo: $q_k^s \rightarrow 0$. Lo contrario también es válido: si $\alpha_k < 1$, un nodo se verá relativamente sobreexpresado. Podemos notar que $\alpha_k < 1$ implica

un valor de expresión $q_k > 1$. Aunque en lógica difusa se supone que los valores de las variables están limitados al intervalo $0 \leq q_k \leq 1$, los valores mayores que uno no están necesariamente excluidos por el formalismo, y es una cuestión de conveniencia el rango en el que se definen las variables [25].

La teoría de redes reguladoras continuas ha sido empleada en el estudio de diferenciación y plasticidad de células del sistema inmune [19, 31], así como para modelar procesos involucrados en el desarrollo celular en plantas [12, 20, 21, 30], o la integración de factores moleculares, celulares y biomecánicos en dichos procesos [32]. Asimismo, este enfoque se utilizó para describir vías de señalización en las células beta del páncreas y la influencia de la inflamación y la sobrenutrición en el desarrollo de la diabetes tipo 2 [?]. En este caso, es posible entender a esta enfermedad como una sucesión de transiciones de fase entre estados de salud, síndrome metabólico y diabetes manifiesta, ocasionadas por el agotamiento de los procesos celulares que permiten compensar la resistencia a la acción de la insulina por parte de estas células.

Cabe hacer notar que el marco conceptual y metodologías involucradas en el estudio de redes neuronales [26, 27] tienen una estructura muy similar a las empleadas en el formalismo de redes difusas; en particular, las interacciones sinápticas se describen mediante funciones equivalentes a la función característica $\mu[w_k]$; es factible suponer que patrones de comportamiento generales de una teoría pueden verse reflejados en la otra.

A. Autoorganización y orden temporal

El desarrollo de estructuras organizadas y asociación con transiciones de fase fuera del equilibrio contempla el hecho de que la transición de un estado inicial a otro final es un proceso esencialmente irreversible, de modo que los cambios entre diferentes patrones emergentes en redes complejas exhiben direccionalidad temporal [12, 21, 28]. Como vimos anteriormente, alteraciones en parámetros del sistema inducen modificaciones del espacio de configuraciones capaces de generar transiciones entre diferentes estados de organización. En un buen número de sistemas complejos un papel central lo juegan las tasas de decaimiento involucradas en la dinámica. Estos parámetros determinan tiempos característicos de expresión de los elementos de la red, definidos por $\tau_k = 1/\alpha_k$. Si suponemos que se puede construir un ordenamiento $\alpha_1 > \alpha_2 > \dots$, es claro que este procedimiento induce un ordenamiento asociado $\tau_1 < \tau_2 < \dots$. En consecuencia, el paisaje de estados puede explorarse modificando las alturas relativas de sus colinas y valles mediante la variación de los parámetros de control α_k , en otras palabras, alterando los tiempos relativos de expresión de las componentes del sistema [30]. Una consecuencia inmediata es que la dinámica del sistema puede discriminar entre variables lentas, $q_L(t)$, y rápidas, $q_R(t)$. Se puede demostrar que la conducta general del sistema a largo plazo está determinado por las variables lentas, mientras que las rápidas se adaptan en forma casi instantánea al entorno definido por las primeras. Lo anterior se refleja en la condición $dq_R/dt \approx 0$, que se cumple para tiempos car-

acterísticos $\tau_R \ll \tau_L$, de modo que q_R alcanza con rapidez su valor estacionario q_R^s . Esta propiedad permite introducir la hipótesis adiabática [2], la cual consiste en la eliminación de las variables rápidas introduciendo su valor estacionario q_R^s en las ecuaciones dinámicas. Por lo tanto, el número de grados de libertad efectivos se reduce y eventualmente limitarse a un número pequeño de variables independientes relevantes, q_I , las pueden identificarse como parámetros de orden del sistema. Este criterio debe tomarse como un grano de sal, dado que cada sistema reviste propiedades dinámicas específicas. En particular, en la teoría de redes las variables que tienen un alto grado de conectividad, así como aquellas que forman 'switches' de bifurcación son centrales en la inducción de patrones autoorganizados.

La existencia de jerarquías de valores de los parámetros característicos de expresión de los agentes que definen un sistema complejo permite establecer una la flecha temporal en la emergencia secuencial de sus diferentes estados de organización [21, 28]. Este mecanismo ha sido empleado para describir la sucesión temporal en el surgimiento de los sépalos, pétalos, estambres y carpelos en las flores [12, 21, 30]. También ha sido utilizado para caracterizar la progresión a largo plazo de la diabetes tipo 2, en donde se observan transiciones graduales entre estados de salud, síndrome metabólico y diabetes manifiesta.

V. REDES ESTOCÁSTICAS

Un elemento relevante en el estudio de las transiciones entre estados estacionarios es la existencia del ruido, es decir, la influencia de interacciones aleatorias inherentes a cada sistema biológico. Dependiendo de su intensidad, la existencia de ruido puede alterar drásticamente las predicciones producidas por el formalismo determinista, especialmente en los puntos de bifurcación del paisaje, donde el ruido puede acelerar la velocidad de transición entre los atractores vecinos [33]. La acción del ruido puede incorporarse en el formalismo de redes reguladoras continuas suponiendo que éste queda representado por una variable estocástica, $\xi(t)$, con promedio nulo $\langle \xi(t) \rangle = 0$, y dispersión estadística dada por $\langle \xi(t) \xi(t') \rangle = DG(t-t')$. Aquí, D es un coeficiente de difusión, y $G\delta(t-t')$ es una función que caracteriza la duración de la autocorrelación de la variable ξ [33]. El caso en el que la correlación es proporcional a la función delta de Dirac, $G(t-t') = G_0\delta(t-t')$, corresponde a un ruido blanco sin efectos de memoria. En el caso más general, la función $G(t-t')$ proporciona una medida del tiempo característico en que el sistema guarda memoria de su estado en tiempos previos, es decir, se trata de un sistema no-markoviano.

El análisis de redes reguladoras con perturbaciones aleatorias ha sido implementado utilizando diferentes tipos de aproximaciones. En la Ref. [12] se contempla un esquema booleano sujeto a fluctuaciones estocásticas capaces de impulsar transiciones entre los valles del paisaje epigenético. En el mismo trabajo se propone un esquema híbrido basado en la dinámica continua por intervalos de Glass [6]. Otros enfoques [21, 34] han considerado la ecuación estocástica de Langevin

[2, 33], la cual tiene la misma estructura que la ecuación determinista (7), excepto por la introducción de un ruido aditivo $\xi(t)$:

$$\begin{aligned} \frac{dq_k}{dt} &= \mu [w_k(q_1, \dots, q_n)] - \alpha_k q_k + \xi_k(t) \\ &\equiv F_k(\mathbf{q}) + \xi_k(t), \end{aligned} \quad (9)$$

con $\mathbf{q} = (q_1, \dots, q_n)$, y en donde $F_k(\mathbf{q})$ juega el papel de una fuerza efectiva. En este esquema, un agente del sistema desarrollará una trayectoria en el espacio de configuraciones guiada por $F_k(\mathbf{q})$, pero sujeta a las fluctuaciones inducidas por $\xi(t)$. Si suponemos que el agente parte del mismo valor inicial $q_k(0)$ en diferentes realizaciones del proceso, esto permitirá construir un ensamble estadístico de trayectorias asociadas a dicha condición. En el caso de ruido pequeño ($D \ll 1$), las soluciones dependientes del tiempo estarán compuestas por la trayectoria promediada sobre el ensamble $\langle q_k(t) \rangle$, más las fluctuaciones aleatorias alrededor de esta ruta. Entonces, los estados estacionarios del sistema están determinados por la relación: $\langle q_k^s \rangle = \langle \mu [w_k^s] \rangle / \alpha_k$. Concluimos que, similarmente al enfoque determinista, los parámetros de decaimiento α_k controlan las alturas relativas de picos y valles en el paisaje epigenético medio. Sin embargo, como se mencionó anteriormente, las fluctuaciones pueden facilitar el tránsito entre valles del paisaje.

Un enfoque complementario es el formalismo de Fokker-Planck [2, 33], el cual considera que el ensamble de trayectorias estocásticas $\{q_k(t)\}$ puede ser caracterizado por su distribución de probabilidad $P[\mathbf{q}(t)]$. La distribución $P[\mathbf{q}(t)]$ satisface la ecuación de difusión:

$$\frac{\partial P[\mathbf{q}]}{\partial t} = - \sum_k \frac{\partial F_k(\mathbf{q}) P[\mathbf{q}]}{\partial q_k} + \frac{1}{2} \sum_{i,k} D_{ik} \frac{\partial^2 P[\mathbf{q}]}{\partial q_i \partial q_k}, \quad (10)$$

en donde $F_k(\mathbf{q})$ es la misma la fuerza efectiva que aparece

en la Ec.(9), y la matriz D_{ik} representa los coeficientes de difusión. Cuando $D_{ik} = D\delta_{ik}$ la difusión de la probabilidad ocurre con la misma tasa D para todas las componentes del sistema. Es fácil mostrar que en el caso en que $F_k(\mathbf{q})$ es derivable de un potencial, $F_k(\mathbf{q}) = -\partial V(\mathbf{q})/\partial q_k$, entonces la ecuación de Fokker-Planck (10) tiene una solución estacionaria dada por:

$$P[\mathbf{q}] = N e^{-2V(\mathbf{q})/D}. \quad (11)$$

Vemos que la estructura de esta solución tiene interpretación muy intuitiva, dado a cada mínimo del potencial $V(\mathbf{q})$ le corresponde un máximo de la distribución de probabilidad $P[\mathbf{q}]$, lo que sugiere proponer el concepto de paisaje de probabilístico [21, 35]. En este caso, los picos más altos de probabilidad de expresión se asientan sobre las cuencas de atracción más profundas del paisaje configuracional, mientras que los valles de probabilidad se localizan en sus cimas. En estos términos, la evolución temporal del sistema puede visualizarse como un tránsito dinámico de los picos de probabilidad en el paisaje de configuraciones. Las redes complejas estocásticas ha permitido entender el orden temporal del surgimiento de los sépalos, pétalos, estambres y carpelos en las flores de la *Arabidopsis thaliana* [12, 21, 30]. Por otro lado, actualmente están en marcha investigaciones enfocadas en la descripción del deterioro paulatino de la funcionalidad celular en las enfermedades crónicas.

A modo de conclusión, el uso de redes reguladoras, ya sea en el marco discreto de Boole, en el continuo involucrado por la lógica difusa y sus generalizaciones estocásticas de Langevin y Fokker-Planck nos puede conducir a una mayor comprensión de los mecanismos subyacentes en la evolución de la complejidad de los sistemas vivos.

-
- [1] W. C. Willett, *Science* **296** 695 (2002)
- [2] H. Haken, *Rev. Mod. Phys.* **47**, 67 (1975); *Synergetics. An Introduction Nonequilibrium Phase Transitions and Self-Organization in Physics, Chemistry and Biology*, Springer-Verlag, Berlin (1978).
- [3] J.D. Estes, S. W. Wietgreffe, T. Schacker, P. Southern, G. Beilman *et al.*, *J. Infect. Diseases*, **195**, 551 (2007). M. Zeng, A.J. Smith, S. W. Wietgreffe, P. Southern, T. Schacker, *et al.*, *J. Clin. Invest.* **121**, 998 (2011).
- [4] Waddington CH. *The Strategy of the Genes*. London: George Allen and Unwin (1957).
- [5] S. Kauffman, *J. Theor. Biol.* **22**, 437 (1969).
- [6] L. Glass, S.A. Kauffman, *J. Theor. Biol.* **39**, 103 (1973); L. Glass, *J. Theor. Biol.* **54**, 85 (1975).
- [7] R.S. Wang, A. Saadatpour, R. Albert, *Phys. Biol.*, **9**, 055001 (2012); R. Albert, J.Thakar, *WIREs Syst. Biol. Med.* **6**, 353 (2014).
- [8] *Gene Regulatory Network Models for Floral Organ Determination*, E. Azpeitia, J. Dávila-Velderrain, C. Villarreal, E. R. Alvarez-Buylla, *Methods in Molecular Biology*, Ed. by JL Riechmann and F. Wellmer, Springer, NY (2014)
- [9] J. Enciso, R. Pelayo, C. Villarreal, *Front. in Immun.*, doi.org/10.3389/fimmu.2019.01927 (2019).
- [10] M. Aldana, *Proc. Nat. Acad. Sci.*, **100**, 8710 (2003); M. Aldana, *Phys. D Nonlinear Phen.*, **185**, 45 (2003); R. Albert, *J. Cell. Sci.*, **118**, **118**, 4947 (2005)
- [11] C. Espinosa-Soto C, P. Padilla-Longoria P, E. R. Alvarez-Buylla, *Plant Cell* **16**, 2923 (2004).
- [12] E.R. Alvarez-Buylla, A. Chaos, M. Aldana, M. Benítez, Y. Cortés-Poza, C. Espinosa-Soto, *et al.*, *PLoS ONE*, **3**, 1 (2008).
- [13] E. Azpeitia, M. Benítez, I. Vega, C. Villarreal, E. Alvarez-Buylla, *BMC Syst. Biol.* **4**, 134 (2010)
- [14] R. Barrio, A. Hernández-Machado, C. Varea, J. R. Romero-Arias, E. R. Alvarez-Buylla, *PlosOne*, **5**, e13523 (2010); E. R. Alvarez-Buylla, E. Azpeitia, R. Barrio, M. Benítez, P. Padilla-Longoria, *Semin. Cell. Dev. Biology* **21**, 108 (2010)
- [15] E. Ortiz-Gutiérrez, K. García-Cruz, E. Azpeitia, A. Castillo, M. P. Sánchez, and E. R. Alvarez-Buylla, *PLOS Comput. Biol.* **11**, e1004486 (2015).
- [16] J. Enciso, H. Mayani, L. Mendoza, R. Pelayo, *Front. in Physiol.* **7**, 349 (2016).
- [17] M.E. Martínez-Sánchez, L. Mendoza, C. Villarreal, E.R.

- Alvarez-Buylla, PLoS Comp. Biol. **11**, 1 (2015).
- [18] D. Martínez-Méndez, C. Villarreal, L. Mendoza, L. Huerta, Front. Immun. (2020). doi.org/10.3389/fphys.2020.00380, (2020)
- [19] L.A. Mendoza, Biosystems **84**, 101 (2006); P. Martínez-Sosa and L.A. Mendoza, Bio Syst **113**, 96 (2013); L.A. Mendoza, BullMath Biol. **75**, 1012 (2013); L.A. Mendoza and A. Méndez, Biosystems. **137**, 26 (2015); A. Méndez and L.A. Mendoza, PLOS Comput. Biol. **12**, e10004696, (2016).
- [20] Y. E. Sánchez-Corrales, E. R. Alvarez-Buylla, L. Mendoza, J. Theor. Biol., **264**, 971 (2010).
- [21] C. Villarreal, P. Padilla-Longoria and E.R. Alvarez-Buylla, Phys. Rev. Lett. **109**, 1 (2012).
- [22] L.A. Zadeh, Inf. and Control, **8**, 228 (1975).
- [23] B. Kosko, Int. J. Gen. Syst. **17**, 211 (1990).
- [24] D. Dubois, S. Moral, Prade, J. Math. Anal. Appl. **205**, 359 (1997)
- [25] V. Novak, I. Perfilieva, J. Mockor, *Mathematical Principles of Fuzzy Logic*, Springer US, Czech Republic (1999).
- [26] J.J. Hopfield, Proc Natl Acad Sci USA, **79**, 2554 (1982).
- [27] Ross TJ. Fuzzy Logic With Engineering Applications. West Sussex: John Wiley (2004).
- [28] J. Wang, L. Xu, and E. Wang, Proc Natl Acad Sci USA. **105**, 12271 (2008); J. Wang J, L. Xu, E. Wang, S. Huang, Biophys J. **99**, 29 (2010)
- [29] J. Bronowski, Leonardo, **18**, 254 (1985)
- [30] J. Dávila-Velderrain, C. Villarreal, E. R. Alvarez-Buylla, BMC Syst. Biol. **9**, 20 (2015).
- [31] M.E. Martínez-Sánchez, L. Huerta, E.R. Alvarez-Buylla, C. Villarreal, Front. in Physiol. **9** (877), 1 (2018).
- [32] V. Hernández-Hernández, R.A. Barrio, M. Benítez, N. Nakayama, J.R. Romero-Arias, C. Villarreal, Phys. Biol. **15** (036002), 1 (2018).
- [33] N. G. van Kampen, *Stochastic Processes in Physics and Chemistry*. North-Holland, Amsterdam (1992).
- [34] J. X. Zhou, L. Bruschi, S. Huang, PLoS ONE 6:e14752. doi: 10.1371/journal.pone.0014752 (2011)
- [35] J. Wang, K. Zhang, L. Xu, and E. Wang, Proc. Nat. Acad. Sci. **108**, 8257 (2011).