

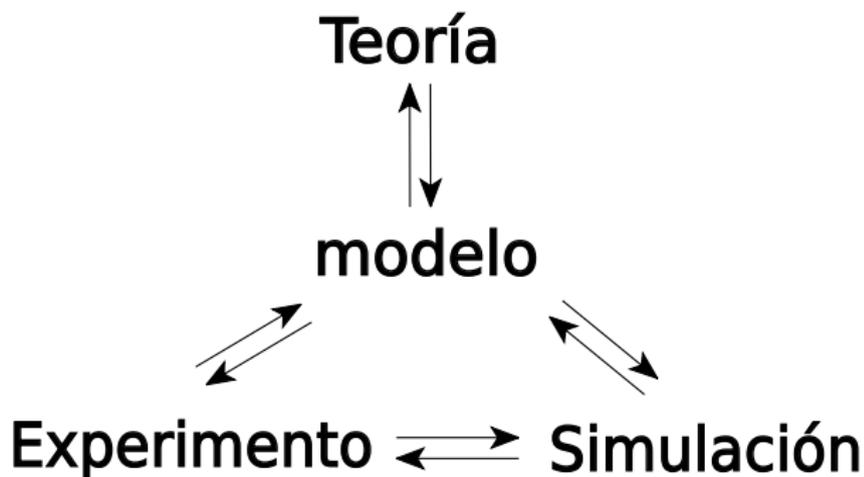
Dinámicas Híbridas QM/MM aceleradas en GPUs

Mariano C. González Lebrero, Matías A. Nitsche y Darío A.
Estrin

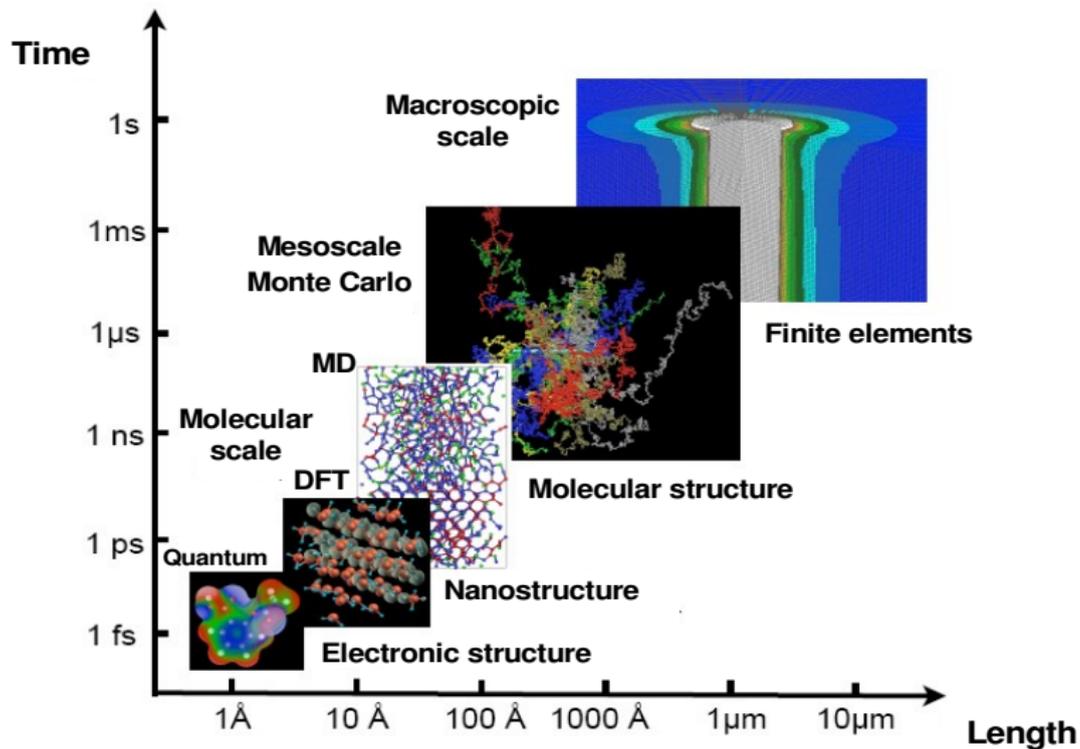
IQUIFIB-CONICET, INQUIMAE-CONICET, FCEN-UBA.

1 de agosto de 2012

Definición: Simular es experimentar con un modelo.



Modelos.



MM

Ya lo vieron ayer...

- Sin detalle electrónico
- Potenciales parametrizados.
- Aplicable a cientos de miles de átomos y en escalas de los nanosegundos o más.

$$\hat{H}\psi = E\psi$$

Parece simple, pero...

- No es razonable, intuitiva ni fácilmente comprensible.
- No se puede resolver exactamente para más de 1 electrón.
- Soluciones autoconsistentes mediante iteraciones.
- Soluciones aproximadas, aunque pueden ser muy buenas.
- Las soluciones mejores sólo se pueden aplicar a unos pocos átomos en el vacío, para todo lo demás existe...

Métodos QM

Hartree Fock

- Semiempíricos.
- HF
- Post Hartree Fock (MP2, CC, CI, etc.)

DFT

- Semiempíricos
- Local Density Aproximation (LDA)
- Generalized Gradient Aproximation (GGAs: PBE, BP86, etc.)

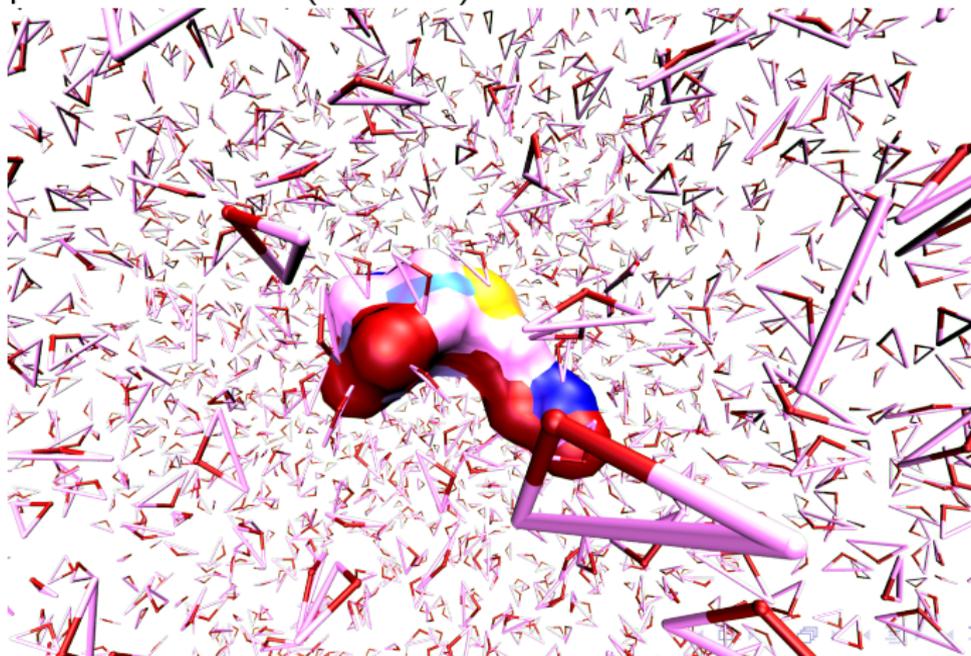
Híbridos (mezcla de DFT y HF)

B3LYP, PBE0, etc.

QM-MM

Métodos Híbridos y Dinámica Molecular

- Para sistemas complejos:
- Métodos híbridos o QM-MM: porción clásica (soluto) + porción cuántica (solvente)



QM/MM y dinámica molecular

$$E = E_{QM} + E_{MM} + E_{QM/MM}$$

$$E_{QM/MM} = \sum_{i=1}^{Ncargas} q_i \int \frac{\rho(r)}{|R_i - r|} dr + \sum_{i=1}^{Clas\ cuan} \sum_{j=1} [V_{LJ}(R) + \frac{q_i Z_j}{|R_i - r_j|}]$$

Para hacer dinámica necesitamos las fuerzas, o sea las derivadas de la energía.

$$F = -\nabla E$$

Haciendo un DFT eficiente...

Partimos de un programa existente, escrito en el año 1992 por Darío Estrin, en fortran 77.

- ➊ **Nuevos algoritmos**, basados en el trabajo de Stratmann *et al.*
 - Menor complejidad
 - Variaciones novedosas, teniendo en cuenta: Procesador gráficos y los tipos de sistemas moleculares que se busca tratar
- ➋ **Nueva implementación**
 - Extiende el programa existente (Fortran77)
 - Se modifica la porción más significativa del cálculo (XC) para hacerlo en la GPU.
 - Código para procesador gráfico
 - Particularidades del hardware y evaluación de alternativas de implementación
 - Optimizaciones a mansalva (mayormente reordenamientos de loops) y uso de memoria dinámica.
 - Empleo de bibliotecas MKL de intel (parelelas).
 - otros.

Método de Kohn-Sham

La energía

$$E[\rho] = T_s[\rho] + V_{ne}[\rho] + \frac{1}{2} \int \int \frac{\rho(\vec{r}_1)\rho(\vec{r}_2)}{r_{12}} d\vec{r}_1 d\vec{r}_2 + E_{xc}[\rho]$$

E_{xc} : energía de intercambio y correlación

$$\rho(\vec{r}) = \sum_{i=1}^{N_{oc}} |\chi_i(\vec{r})|^2$$

$$\chi_i(\vec{r}) = \sum_{\mu=1}^M C_{i\mu} \phi_k(\vec{r})$$

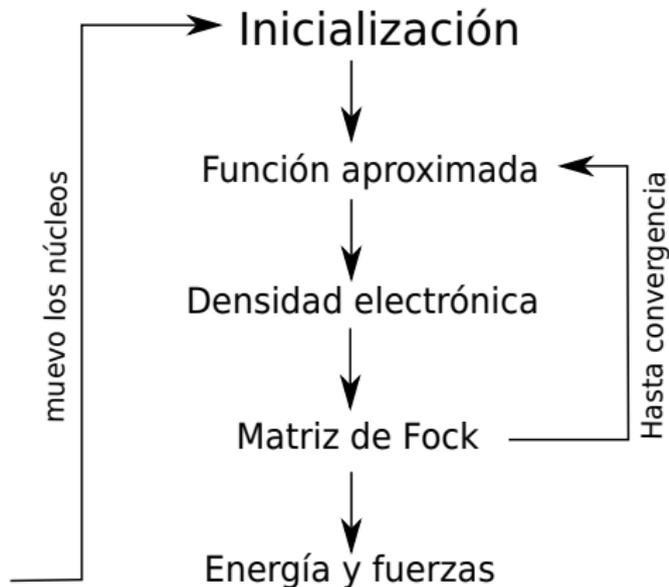
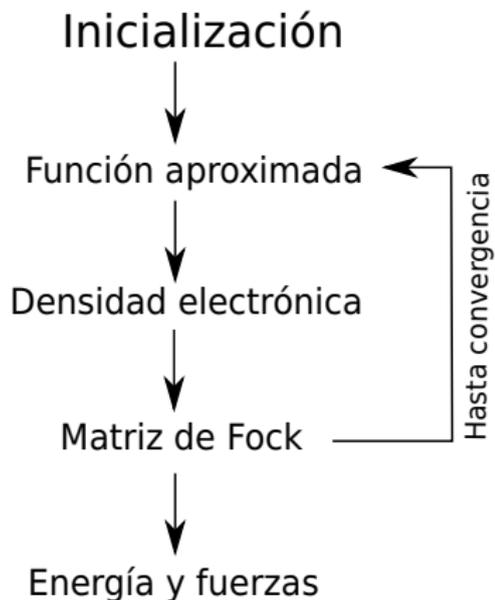
Las ϕ_k son funciones llamadas base, en nuestro caso son Gaussianas.

... se llega a que

$$\rightarrow \rho(\vec{r}) = \sum_{i=1}^{N_{oc}} \left| \sum_{k=1}^M C_{ik} \phi_k(\vec{r}) \right|^2$$

Se trata de encontrar los coeficientes C_{ik} que minimicen la energía.

Esquema de cálculo



Energía de Intercambio y Correlación

- Múltiples métodos de aproximación, el más simple es LDA:

$$E_{xc}[\rho] = \int \rho(\vec{r}) \varepsilon_{xc}(\rho(\vec{r})) d\vec{r}$$

ε_{xc} : funcional de intercambio y correlación

- Si sólo depende de la densidad.
- Otros método más preciso: GGA (Global Gradient Approach).

$$E_{xc}[\rho] = \int \rho(\vec{r}) \varepsilon_{xc}(\rho(\vec{r}), \nabla \rho(\vec{r})) d\vec{r}$$

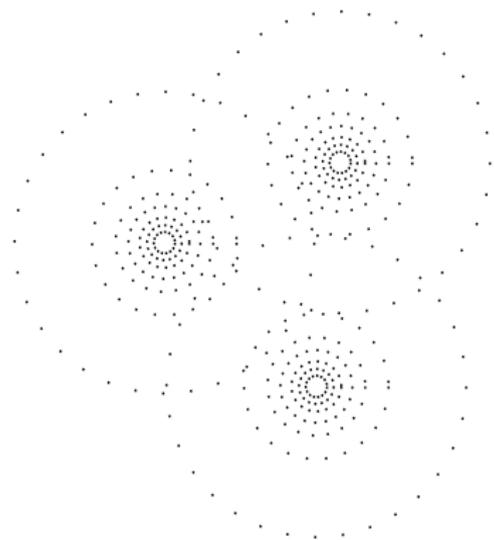
No existe solución analítica de esta integral, se integra numéricamente sobre una grilla.

$$I = \int F(\vec{r}) d\vec{r}$$

$$I \simeq \sum_i \omega_i F(\vec{r}_i)$$

La Grilla

- Se toma una grilla esférica como base.
- Superposición de grillas base, c /escalamiento: capas concéntricas, sobre cada átomo
- En sistemas poliatómicos: unión de las grillas de c /átomo



Las funciones.

La función de onda se describe como combinación lineal de un grupo de funciones, llamado Base.

La base puede ser de funciones Gaussianas, ondas planas, etc. En nuestro caso son Gaussianas, es decir tienen esta pinta:

$$\phi = (x - x_0)^n (y - y_0)^m (z - z_0)^l e^{\alpha |\vec{r} - \vec{r}_0|^2}$$

Donde n , m , y l son número enteros que dan el momento angular a la función.

Usar funciones Gaussianas es venajoso porque reproducen adecuadamente la distribución electrónica y son matemáticamente amigables.

La densidad.

Los electrones están en orbitales moleculares:

$$\chi_i(\vec{r}) = \sum_{\mu=1}^M C_{i\mu} \phi_k(\vec{r})$$

Entonces la densidad electrónica total es la suma de las contribuciones de cada orbital:

$$\rightarrow \rho(\vec{r}) = \sum_{i=1}^{N_{oc}} \left| \sum_{k=1}^M C_{ik} \phi_k(\vec{r}) \right|^2$$

O dando vuelta la sumatoria se puede llegar a que:

$$\rightarrow \rho(\vec{r}) = \sum_{j=1}^M \sum_{k=1}^M P_{jk} \phi_k(\vec{r}) \phi_j(\vec{r})$$

Donde P_{ij} es la llamada matriz densidad.

Empezando a pensar algoritmos eficientes.

Costo de calcular la integral de intercambio y correlación = $\alpha M^2 \times N$ donde M es el número de funciones base y N el número de puntos de la grilla.

¿Que hacemos?

Respuesta: funciones significativas.

funciones significativas.

- 1 Agrupar los puntos de la grilla.
- 2 Encontrar las funciones significativas
 ϕ es significativa en un grupo si $\phi(r_i) \geq \epsilon$ para algún punto (o lugar) del grupo.

Ahora hay que definir el tamaño de los grupos y el valor de ϵ .

- Grupos muy chicos \rightarrow menos funciones x grupo pero más grupos.
- ϵ grande \rightarrow menos funciones significativas pero mayor error.

Si hacemos esto para sistemas grandes el número de funciones significativas no crece con el tamaño del mismo.

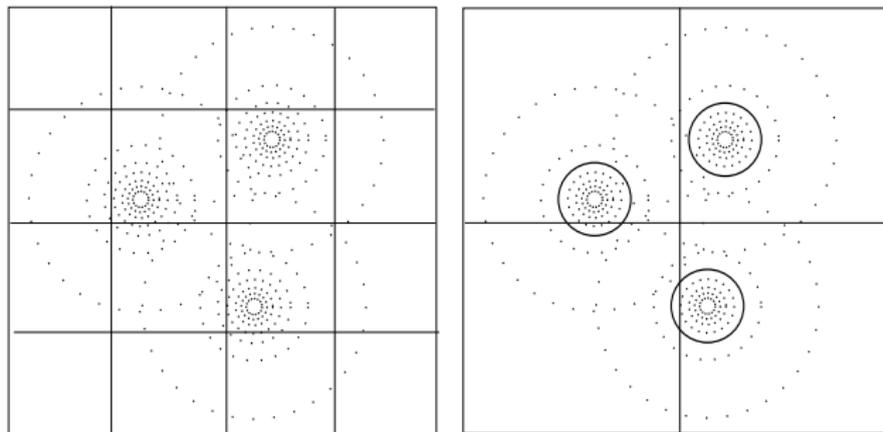
¡Escalamiento lineal!

Algoritmos eficientes.

Ej: Particionado del sistema y funciones significativas.

Algoritmo usual basado en cubos \Rightarrow Muchos cubos con pocos puntos (malo para GPU).

Algoritmo basado en esferas y cubos \Rightarrow Mucho más eficiente.



El nuevo elefante: Coulomb

$$E_{Coul} = \frac{1}{2} \int \int \frac{\rho(\vec{r}_1)\rho(\vec{r}_2)}{r_{12}} d\vec{r}_1 d\vec{r}_2$$

Donde

$$\rho(r) = \sum_{i(ocup)} \sum_{j(Base)} \sum_{k(base)} C_{ij} C_{ik} \Phi_i \Phi_j = \sum_{i(base)} \sum_{j(base)} P_{ij} \Phi_i \Phi_j$$

$$\rho \cong \tilde{\rho} = \sum_{Auxiliar} C_k G_k$$

$$E_{Coul} \approx \int \int \frac{\rho(\vec{r}_1)\tilde{\rho}(\vec{r}_2)}{r_{12}} d\vec{r}_1 d\vec{r}_2 - \frac{1}{2} \int \int \frac{\tilde{\rho}(\vec{r}_1)\tilde{\rho}(\vec{r}_2)}{r_{12}} d\vec{r}_1 d\vec{r}_2$$

Coulomb

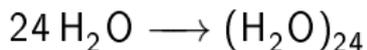
- Se pueden tirar términos basado en el teorema de productos de Gaussianas.

$$\int e^{-a(\vec{r}-\vec{r}_a)^2} \times e^{-b(\vec{r}-\vec{r}_b)^2} = \int e^{-\frac{ab}{a+b}(\vec{r}_a-\vec{r}_b)^2} e^{-(a+b)(\vec{r}-\frac{a\vec{r}_a+b\vec{r}_b}{a+b})^2}$$

- Escala cuadráticamente.
- Se puede guardar en memoria y no recalcular en cada iteración (los productos de las Gaussianas no cambian, sólo los coeficientes P_{ij})
- El cálculo de las derivadas (para la fuerza) es lo más demandante para sistemas medianos (HEMO), junto con la integral QM/MM.

Ootro tema: Precisión.

Energía de formación de un agregado.



$$\nabla E = E((\text{H}_2\text{O})_{24}) - 24E(\text{H}_2\text{O})$$

$$\nabla E = -1832,896H - 24 \times (-76,352H) = 0,437H$$

La precisión química es del orden de 1 kcal/mol o sea 10^{-3} Hartrees!

Es indispensable usar precisión alta siempre, no??

Respuesta: no

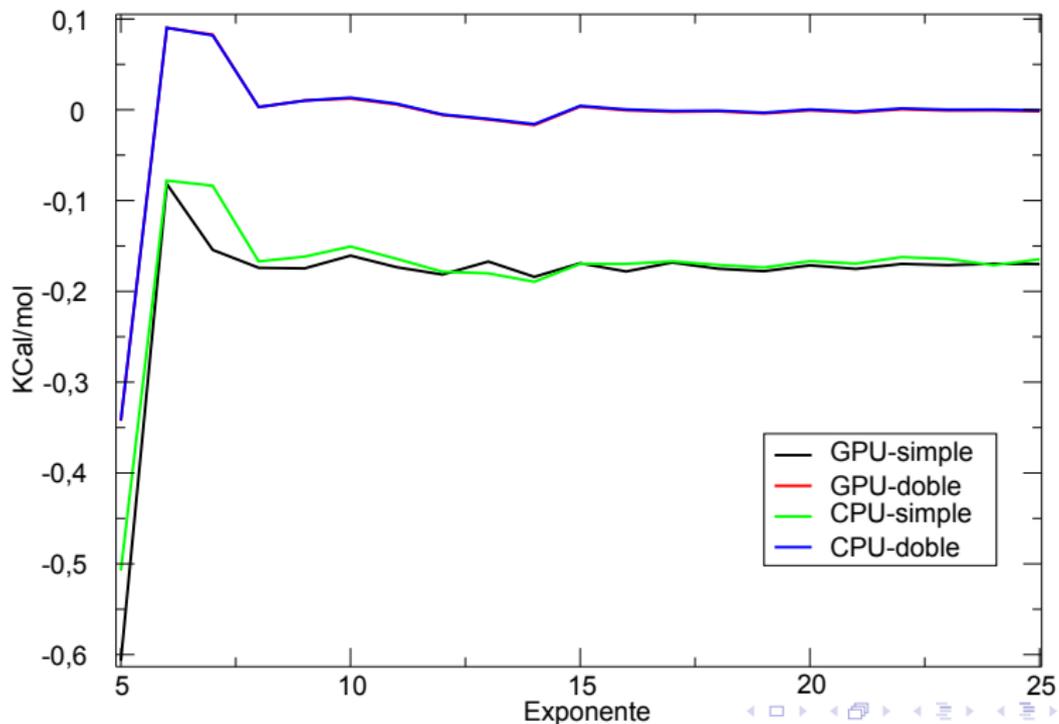
- Se puede usar precisión simple para muchas partes manteniendo la calidad del resultado.
- Para otras se puede usar precisión mixta.
- Mejora los tiempos de cálculo (¡sobre todo en GPUs!).
- Disminuye el uso de memoria (muy útil en Coulomb).

uso de precisión simple.

En la integral de intercambio y correlación para todo salvo para la acumulación.

Para las integrales de Coulomb se aplica el mismo criterio basado en el producto de Gaussianas. Si son muy chicas, se tiran, si tienen un valor intermedio se hacen en simple, si tienen un valor mayor en doble.

Calidad.

Dif. energía de formación
del agua 24

Calidad

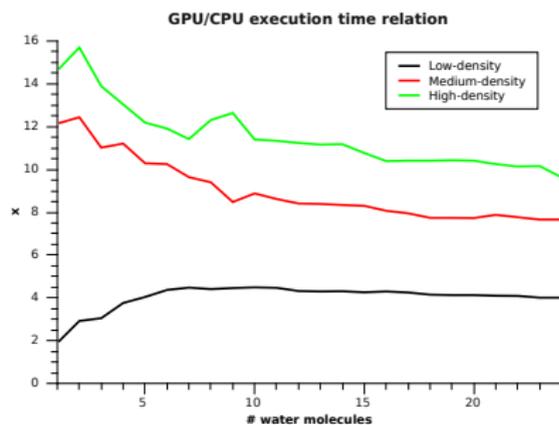
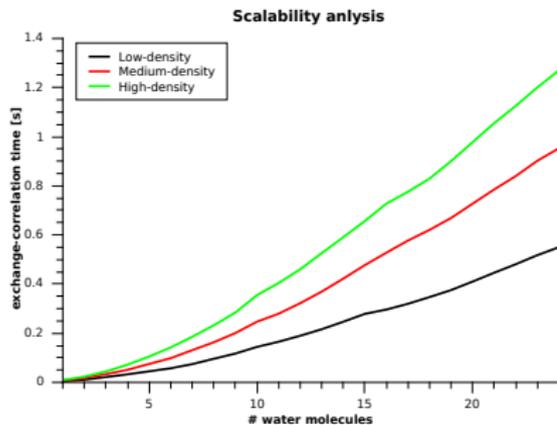
- Usando precision simple se comete un error en la energía de 0,2 Kcal/mol para los sistemas mas grandes (incluso para la valinomicina, de 168 átomos).
- Con exponentes mayores a 10 ya se obtiene un buen resultado.
- GPU en doble es equivalente a CPU en doble.

Parámetros nuevos.

Tienen impacto en la performance y la calidad.

- `max_function_exponent` => Define el valor de las funciones que se tiran.
- Tamaño de los cubos.
- Tamaxo de las esferas (entre 0. y 1.)
- `rmax` => Define los pares de Gaussianas que se tiran cuando hace Coulomb.
- `rmaxs` => Define los pares de Gaussianas que se hacen en doble.

Escalamiento y aceleración.



	Taxol-631G		Valinomicina-631G		Hemo-DZVP	
	GPU	CPU/GPU	GPU	CPU/GPU	GPU	CPU/GPU
Generación de la grilla	0.99	19.53	1.63	18.88	0.42	18.79
Iteración SCF	3.70	8.19	6.40	6.03	2.92	9.13
Intercambio y correlación	2.41	11.96	3.5	10.09	2.3	11.26

Garchamber

Ahora lo bueno, a este código mejorado lo “pegué” al amber.

- Usé buena parte del código de QM/MM de Amber (para DFTB y PM3).
- Necesita un archivo de input extra con las bases y los parámetros del DFT. (y se tiene que llamar input)
- Pase todo a memoria dinámica (digamos un 90 %) y las variables en un módulo (no mas COMM).
- Compilado como biblioteca.
- Cut off, no Ewald.

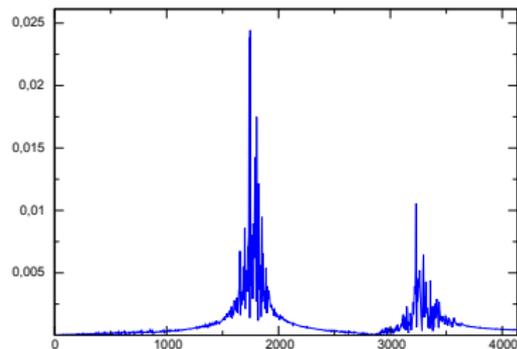
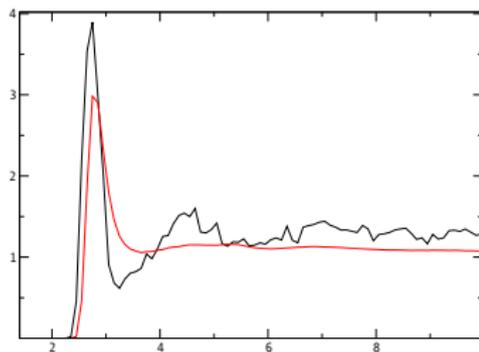
¿Cómo anda?

Sistema	segundos por paso de dinámica
agua1 en agua	0.56
agua2 en agua	1.3
agua3 en agua	2.3
Hemo O2	180
trip en agua	20

Resultados para el agua

40000 pasos de dinámica en 6 horas y media

$g(r)$



¿Qué se logró?

- ① Una implementación de DFT con variaciones novedosas:
 - Esferas → mejora desempeño en GPU
 - En sistemas relativamente pequeños (útiles para QM-MM):
- Mejor rendimiento usando GPU
 - Speedup: entre cerca de $10x$
- Poco uso de memoria
 - Se podría aplicar a sistemas más grandes
- Buena calidad numérica.
 - ¡Solo con precisión simple!
 - Mejor con precisión doble.

Lo aprendido

- Límites/virtudes de los GPU
 - Cómo encarar futuros problemas (*divide & conquer*)
- Aspectos de performance
 - Principal: acceso a memoria
 - Patrón de acceso eficiente, generalizable a otros problemas

Trabajo a futuro (se aceptan voluntarios)

- Hacer mas amigable el input (bases en un formato mas estandar y etc.)
- Hacer UHF (se supone que es fácil)
- Mejorar el método de convergencia.
- Integrales de Coulomb y QM/MM en GPU.
- Outputs de orbitales en un formato graficables.
- ¿Condiciones periódicas de borde?, ¿cutoff por residuo?
- ¿Método lineal para las integrales de coulomb? (para sistemas enormes, Ej: Proteínas enteras.)

Autores.

Lic. Matías Nitsche.

Dr. Darío estrin.

Gracias!