

Voronoi en el GPU

Emmanuel Rojas Fredini

November 24, 2009

Contents

1	Diagrama de Voronoi? se come?	3
1.1	Introduccion	3
1.2	Definicion	3
1.3	Aplicaciones	4
1.4	Calculo	6
2	Introduccion a CUDA	8
2.1	Jerarquia Hilos	8
2.2	Jerarquia de Memoria	8
3	Voronoi + CUDA = Mi Trabajo Final?	10
	List of Figures	11
	List of Tables	12

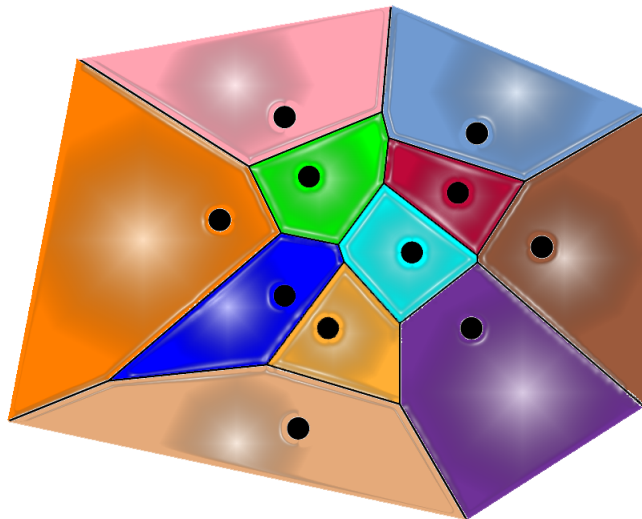
Chapter 1

Diagrama de Voronoi? se come?

1.1 Introduccion

dasgdkshdgsa

1.2 Definicion



References Definamos la distancia Euclidea entre 2 puntos p y q como $dist(p, q)$. Luego en el plano tenemos

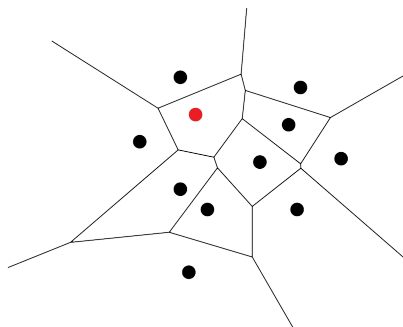
$$dist(p, q) = (p_x - q_x)^2 + (p_y - q_y)^2 \quad (1.1)$$

Y análogamente para mas dimensiones. Sea $P = \{p_1, p_2, \dots, p_n\}$ n puntos distintos en el espacio. Llamaremos a estos puntos sitios. Entonces definimos

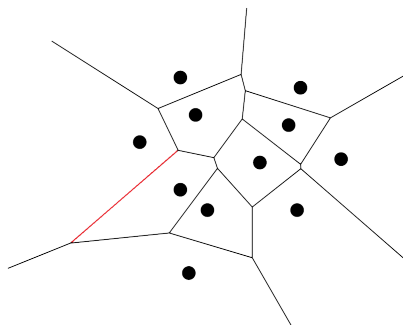
el diagrama de Voronoi de P como la subdivisión del plano(o del espacio) en n celdas, una para cada sitio en P , con la propiedad de que el punto q esta en la celda correspondiente al sitio p_i si $dist(q, p_i) < dist(q, p_j)$ para $p_i \wedge p_j \in P, j \neq i$. Asimismo denotaremos al diagrama de Voronoi de P como $Vor(P)$, y la celda p_i del diagrama anterior la denotamos $Vor(p_i)$.

De esta forma el diagrama de Voronoi esta formado de los siguientes elementos:

- Sitios, puntos en R^2, R^3, \dots, R^n . [?]



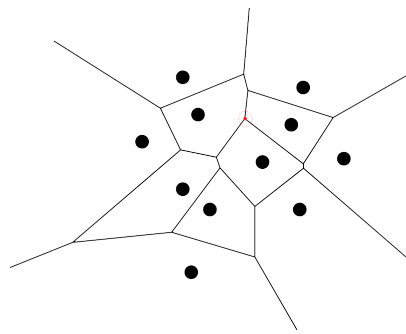
- Aristas o planos de bisección, son los planos(o segmento, hiper-planos) que separan a 2 celdas, delimitando los puntos más cercanos a un sitio p_i y otro p_j con $i \neq j$. Tienen la propiedad de ser equidistantes a 2 sitios. [?]



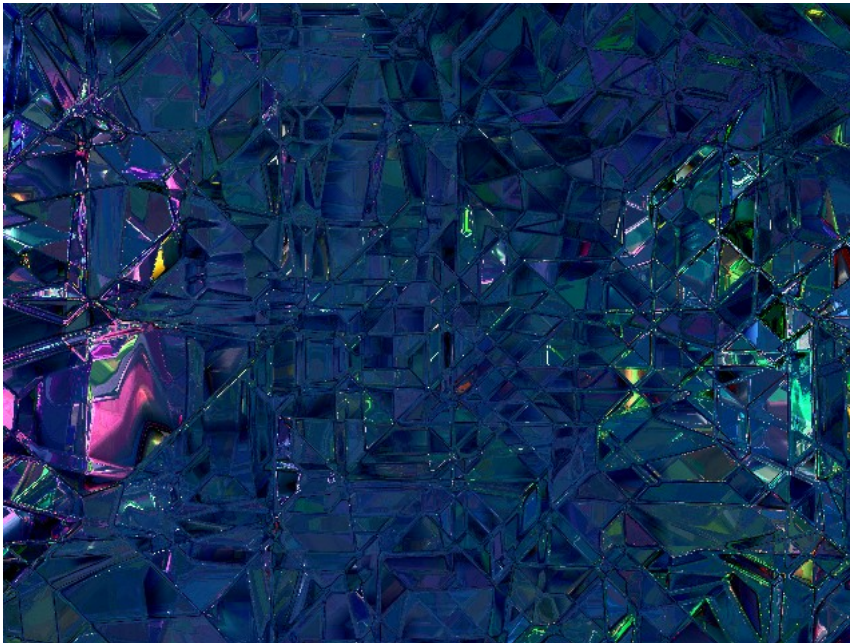
- Semillas o vértices, son las uniones de las aristas o planos de bisección, se caracterizan por que si una circunferencia(o esfera, hiper-esfera) crece desde este chocaría con varios sitios al mismo tiempo, ie son equidistantes a 3 o mas sitios. [?]

1.3 Aplicaciones

Se puede construir una estructura en base al diagrama de Voronoi para encontrar el vecino mas cercano, relacionando a que celda pertenece el punto, y que celdas vecinas tiene, es decir las aristas de la tringulacion de Delaunay. Un uso muy extendido es la cuantización vectorial, muy usado en la compresión de datos.

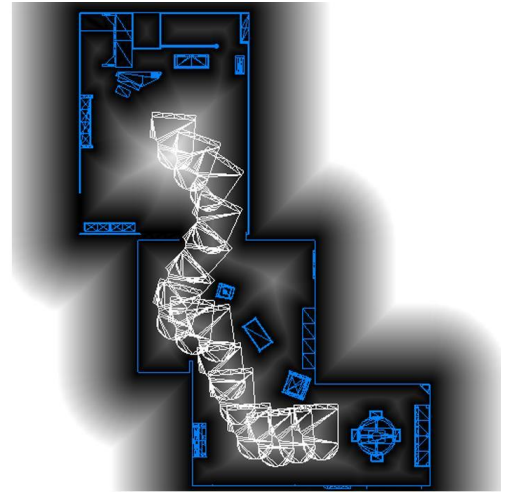
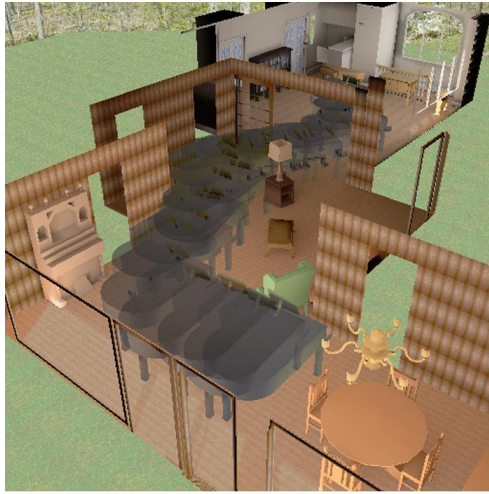


Se puede encontrar el mayor círculo vacío entre todos los sitios, y también polígonos, esto se lo reconoce verificando los vértices o semillas del diagrama. Esta propiedad puede usarse para identificar la posición más conveniente para el establecimiento de un negocio, etc. Es usado en la física de polímeros¹. Con este se puede representar el volumen libre de un polímero. Se usa en el análisis de la capacidad de una red inalámbrica. Se lo utiliza para analizar los patrones de crecimiento de bosques, y en predicciones de la expansión de incendios forestales. También se los usa para generar proceduralmente texturas con apariencia orgánica[?].



En robótica, se usa el diagrama de Voronoi para encontrar rutas de movimiento. Si los puntos son obstáculos, entonces las aristas del diagrama serán las rutas más alejadas de los obstáculos[?].

¹Un polímero es una molécula grande (macromolécula) compuesta por una estructura de unidades repetidas, generalmente conectadas por una unión covalente



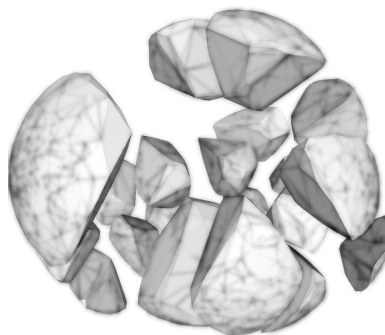
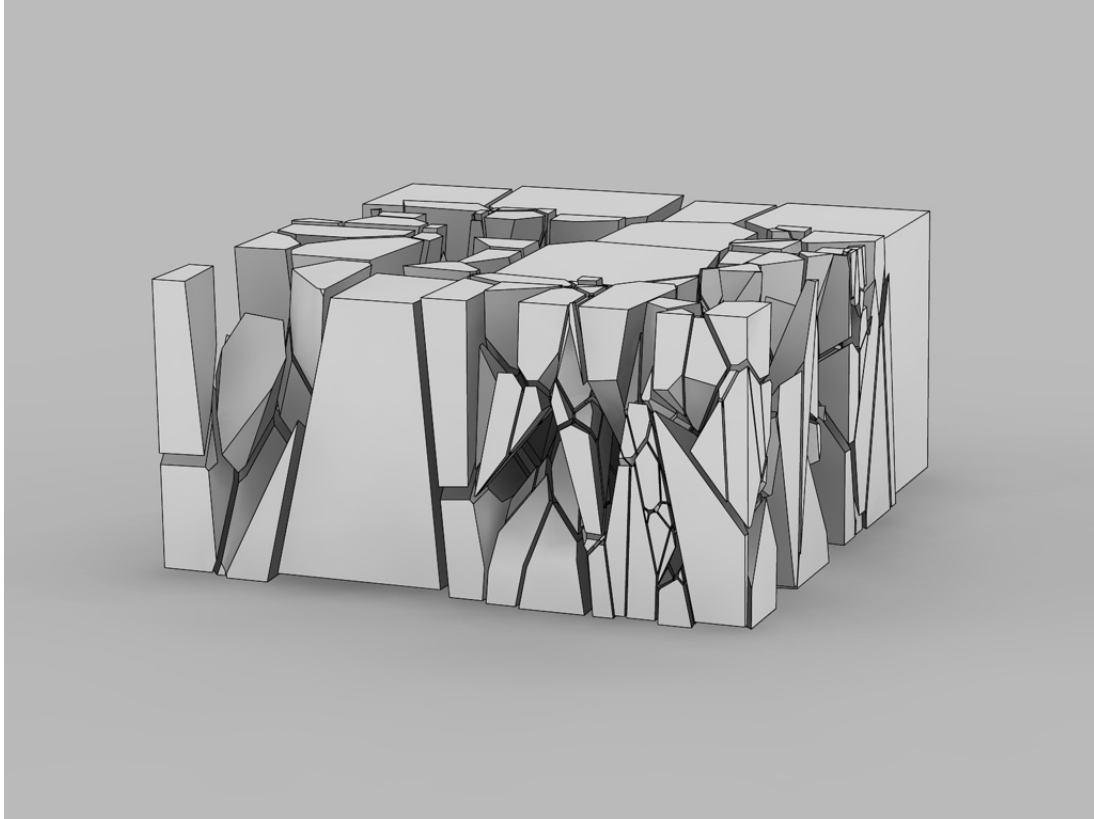
Se usa para teselar modelos 3D en fractales, para luego simular ruptura de elementos. Un ejemplo de esto es el Voronoi Shatter de Blender, que permite teselar un modelo, pudiendo luego hacer una simulacion fisica con el motor fisico de Blender[?][?].

1.4 Calculo

Algunos métodos para calcular este diagrama de forma eficiente son:

1. Linea de Fortune que presenta un orden $O(n * \log(n))$ en el plano y $O(n)$ en el espacio, donde n son los sitios en el espacio.
2. Algoritmo basado en Dividir y Conquistar.
3. Calculo de triangulación de Delaunay, y pasar por la dualidad *Voronoi-Delaunay* al Diagrama de Voronoi
4. Jump Flooding, paper de Guodong Rong y Tiow-Seng Tan, que presenta un orden $O(k)$ donde k es una constante. Este método requiere de un procesador multi nucleo.
5. Renderizado de Voronoi usando funciones de distancia, paper de Hoff, requiere renderizador muli nucleo SIMD.

Explicaremos otro método que utiliza procesadores SIMD (más específicamente una GPU CUDA enabled) para calcular el diagrama de Voronoi de forma discreta, es decir que se calculara en una grilla de axb , o $axbxc$, o generalizado para mas dimensiones, en lugar de en el espacio continuo R^2 , R^3 , o R^n .



shatte2.jpg

Chapter 2

Introduccion a CUDA

CUDA(Compute Unified Device Architecture) es un modelo de programacion paralela con operaciones SIMD(Single Instruction Multiple Data), que permite dar instrucciones a grupos de hilos(threads) de ejecucion mediante un lenguaje compuesto de C con extensiones. Este se basa en una arquitectura de implementacion de procesadores multi-nucleo, que actualmente solo es implementada por Nvidia(su creador). Asimismo existen otros modelos muy similares entre los que destacamos Stream Computing de Ati, OpenCL y ComputeShader sinendo todos estos distintos modelos de programacion SIMT(Single Instuction Multiple Thread), es decir que permiten ejecutar un grupo de instrucciones por gran cantidad de threads en paralelo.

El modelo es inerentemente jerarquico, tanto en su modelo de memoria, como en su modelo de threads. Aumentando la capacidad de granularidad de un algoritmo.

2.1 Jerarquia Hilos

Los hilos se identifican un identificador de 3 variables x , y y z . Esto permite que la paralelisacion en base a arreglos de hasta 3 dimensiones sea facil de implementar, y no se deba recurrir a mapeos oscuros que ofusquen el codigo. A su vez los threads se agrupan en bloques, identificandose estos por 2 variables x e y . También permitiendo esto simplificar los algoritmos.

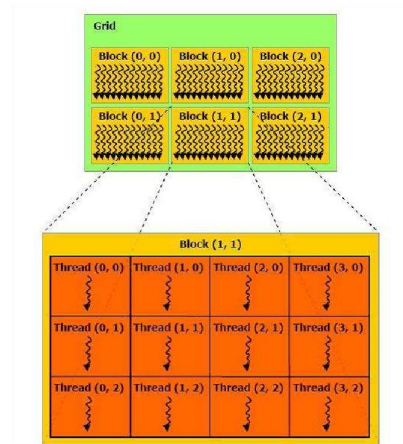
En fin un thread es identificado universalmente por su id de hilo o *threadIdx*, y su id de bloque o *blockIdx*.

Cuando se habla de la grilla(Grid), se refiere al arreglo bi-dimensional de bloques con los que se ejecuta el kernel.

2.2 Jerarquia de Memoria

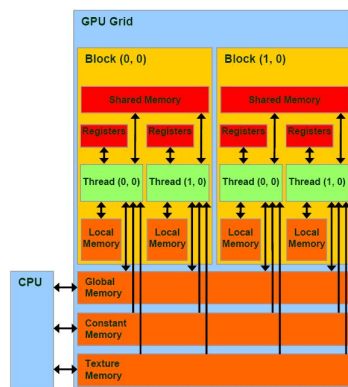
Esta se compone de:

1. Un conjunto de registros pertenecientes a cada thread exclusivamente. Tiene un periodo de vide ligado al thread.



2. Una memoria compartida(Shared Memory) entre todos los threads de un bloque. Tiene un periodo de vida ligado al bloque/thread, ya que estos "nacen y mueren juntos".
3. Una memoria de Textura, es solo de lectura y esta optimizada para mapeos d texturas y sus formatos, compartida por el kernel. Tiene un periodo de vida ligado al kernel.
4. Una memoria Contante, es solo de lectura, compartida por el kernel. Tiene un periodo de vida ligado al kernel.
5. Una memoria global compartida por el kernel. Tiene un periodo de vida ligado al kernel.

La jerarquia fue nombrada en un orden de velocidad de memoria decreciente o identica, y de latencia creciente o identica. Empero en tamaños decedentes o identicos. Luego la memoria compartida, de textura y constante puede considerarse como un cache administrado por software y no por hardware, ya que permite una gran velocidad de acceso y modificacion. Por eso es conveniente copiar la memoria sobre la que se va a trabajar a alguno de estas memorias.



Chapter 3

Voronoi + CUDA = Mi Trabajo Final?

Este metodo de calculo de Voronoi discreto, o aproximacion a Voronoi, utilizando la GPU tiene la ventaja sobre los otros 2 algoritmos para calcular Voronoi en la GPU, el metodo de Hoff(renderizado de funciones de distancia) y Jump Flooding, en que este metodo es facil de modificar para obtener la solucion a otros problemas, como son obtener los puntos mas lejanos y k-esimo Diagrama de Voronoi. Ademas el metodo de Hoff puede tener un cuello de botella cuando se sobrecarga a la GPU con demasiados sitios, es decir que no se escala como el algoritmo que exponaremos.

El algoritmo consiste en ejecutar por cada pixel de la grilla un thread, y en cada thread calcular la distancia al punto mas cercano, luego este pixel pertenece a la celda de ese sitio. Expresando esto en un pseudocodigo:

```
Para cada pixel
{
  Calcular la distancia de este pixel a todo los sitios
  Elegir el sitio mas cercano
  Colorear este pixel con el color del sitio mas cercano
}
```

De esta forma si desea lograr la solucion a los puntos mas lejanos, en vez de buscar la menor distancia, se busca la mayor distancia. Para lograr una buena implementacion con este algoritmo se debe utilizar la memoria compartida para buscar la menor distancia, ademas llamar al kernel con los parametros mas convenientes para la GPU que la ejecuta. Esto se debe a que por ejemplo la placa Nvidia 8800 GTS posee 12 multiprocesadores, de esta forma lo mas conveniente es tener 12 blockes, y para aumentar el uso de cada multiprocesador, que posee 8 procesadores stream y puede ejecutar como maximo 512 threads al mismo tiempo, asi que lo mas conveniente es hacer arreglos de 512 threads.

Otra forma de mejorar el algoritmo es utilizar una estructura de ordenamiento espacial como un kd-tree, un quadtree, o simplemente agrupar por grillas equidimensionadas.

List of Figures

List of Tables