

GPU Computing: una mirada al pasado, presente y futuro

Charla Invitado Universidad del Bio Bio

*Dr. Cristóbal A. Navarro
Instituto de Informatica
Facultad de Ciencias de la Ingeniería
Universidad Austral de Chile*



Introducción



Desafíos

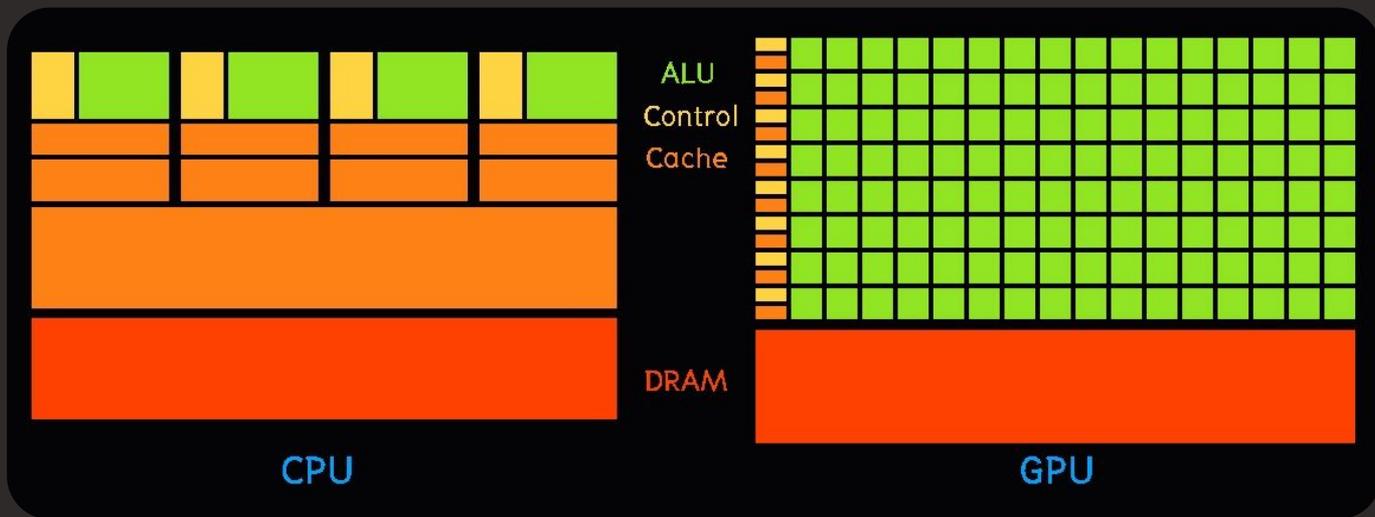


Grupo Temporal

Graphics Processing Unit (GPU)



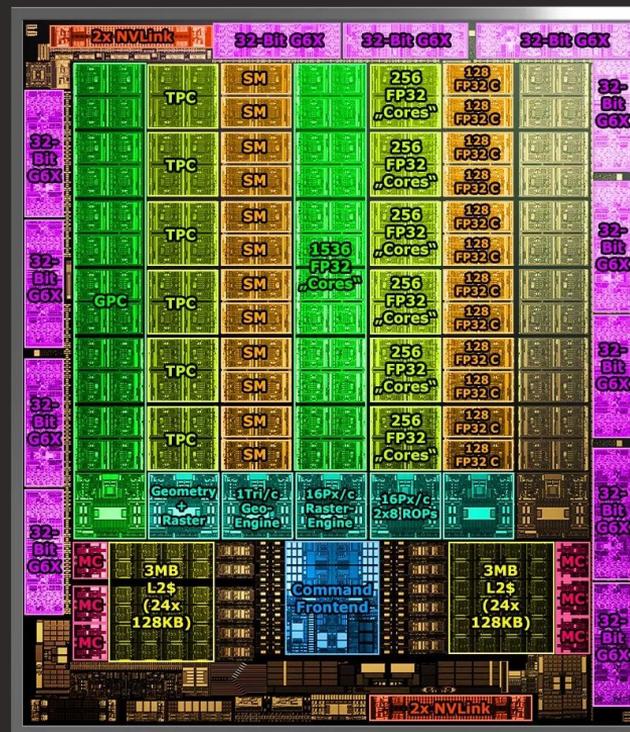
CPU y GPU: Qué las diferencia?



Chips en vista IR



CPU (Alder Lake 2021)



GPU (Ampere, GA102)

Videojuegos cumplen un rol clave



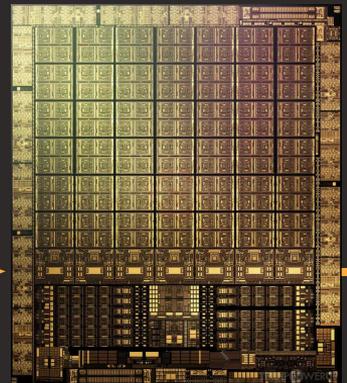
GPU: Arquitectura y Organización



Tarjeta



Chip



Chip (zoom físico)



Chip (zoom lógico)

GPU: Arquitectura y Organización



Algunas GPUs gama alta 2021-2022

NVIDIA A100 SXM4 40 GB

GA100	6912	432	160	40 GB	HBM2e	5120 bit
GRAPHICS PROCESSOR	CORES	TMUS	ROPs	MEMORY SIZE	MEMORY TYPE	BUS WIDTH

~19.5 TFLOPS (FP32)

AMD Radeon RX 6900 XT

Navi 21	5120	320	128	16 GB	GDDR6	256 bit
GRAPHICS PROCESSOR	CORES	TMUS	ROPs	MEMORY SIZE	MEMORY TYPE	BUS WIDTH

~23 TFLOPS (FP32)

Intel Xe-HP (2022)

Scalability

1 Tile 2 Tile 4 Tile

~42 TFLOPS (FP32)

NVIDIA GeForce RTX 3090

GA102	10496	328	112	24 GB	GDDR6X	384 bit
GRAPHICS PROCESSOR	CORES	TMUS	ROPs	MEMORY SIZE	MEMORY TYPE	BUS WIDTH

~36 TFLOPS (FP32)

AMD Radeon Instinct MI100

Arcturus	7680	480	64	32 GB	HBM2	4096 bit
GRAPHICS PROCESSOR	CORES	TMUS	ROPs	MEMORY SIZE	MEMORY TYPE	BUS WIDTH

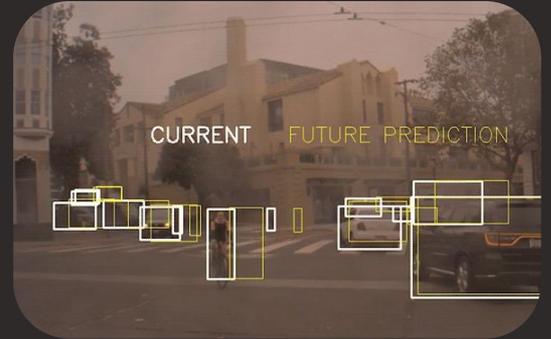
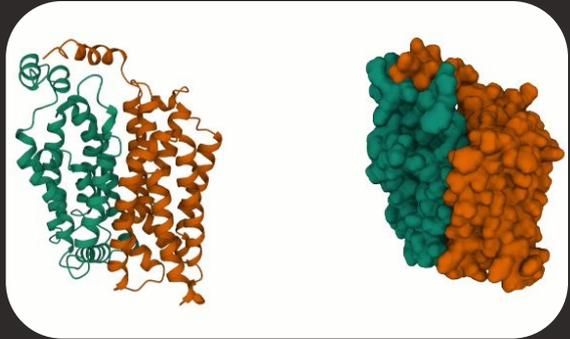
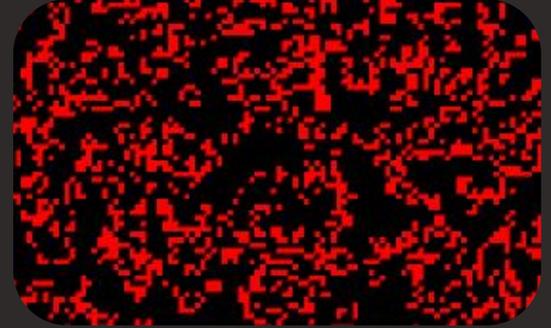
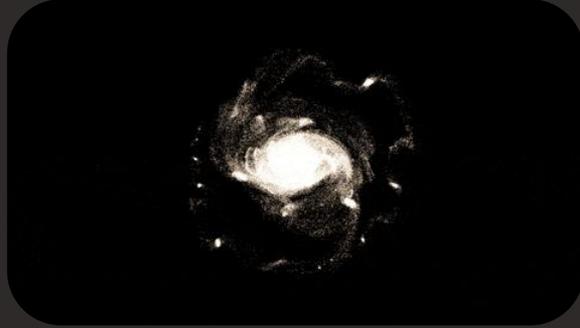
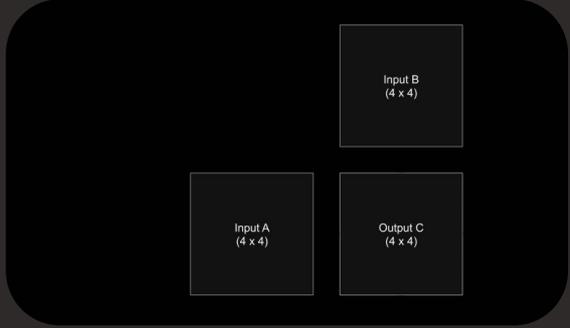
~23 TFLOPS (FP32)

Intel Xe DG1-SDV

DG1	768	48	24	3 GB	LPDDR4X	96 bit
GRAPHICS PROCESSOR	LORES	TMUS	ROPs	MEMORY SIZE	MEMORY TYPE	BUS WIDTH

~2.3 TFLOPS (FP32)

GPU Computing - Aplicaciones



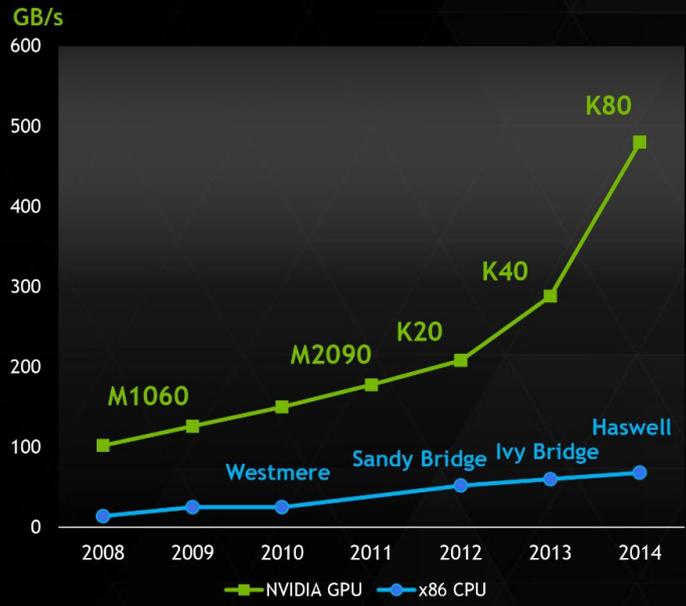


GPU Computing (motivación)

Peak Double Precision FLOPS

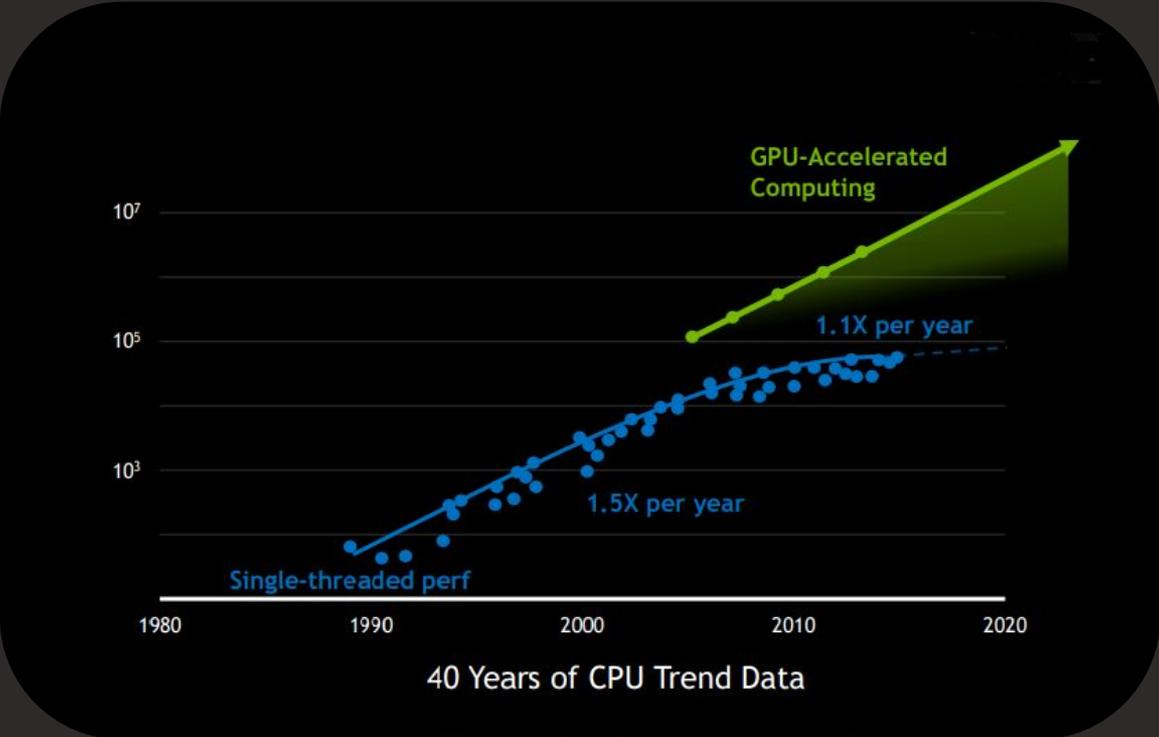


Peak Memory Bandwidth



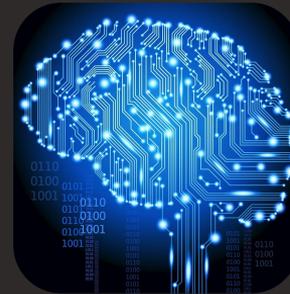
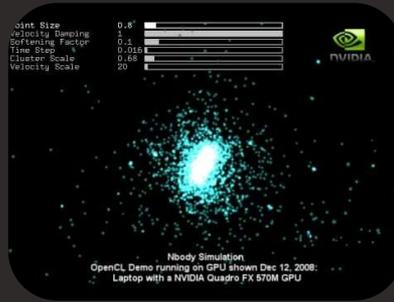


GPU Computing (proyección)



40 Years of CPU Trend Data

Estamos en la 3ª era de GPU Computing



2003

1ª era
(GPU Hacking)

2007

2ª era
(Programacion General GPU)

2017

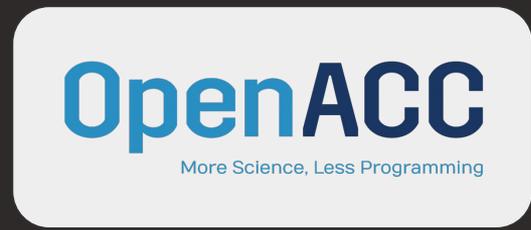
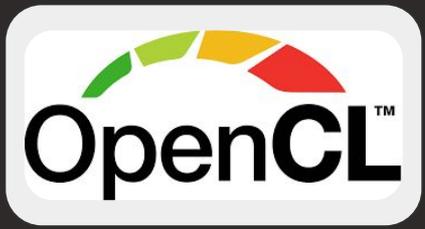
3ª era
(Programacion ASIC)



Programación en GPUs:



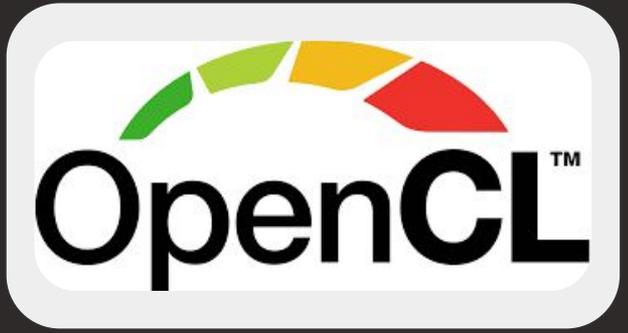
Abstracción específica



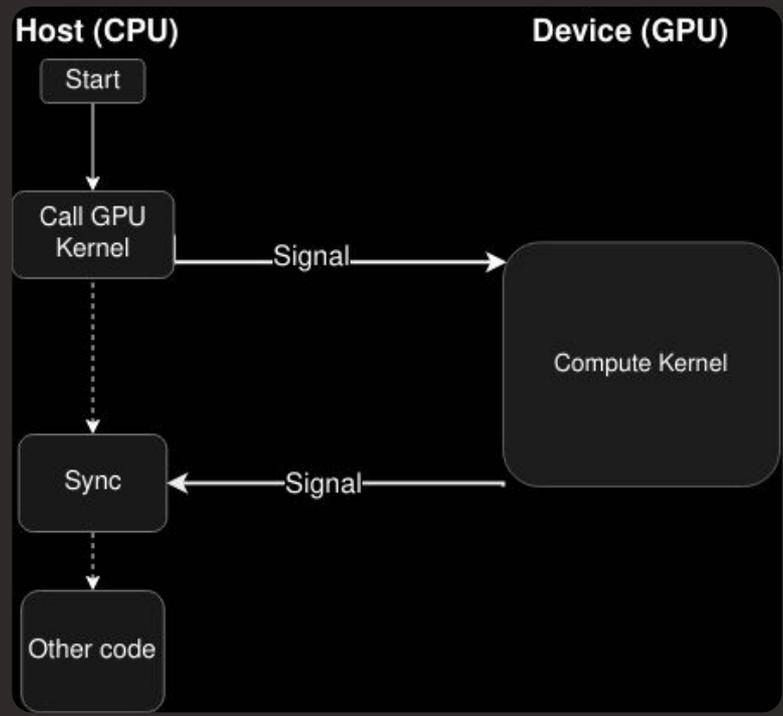
Abstracción HPC General



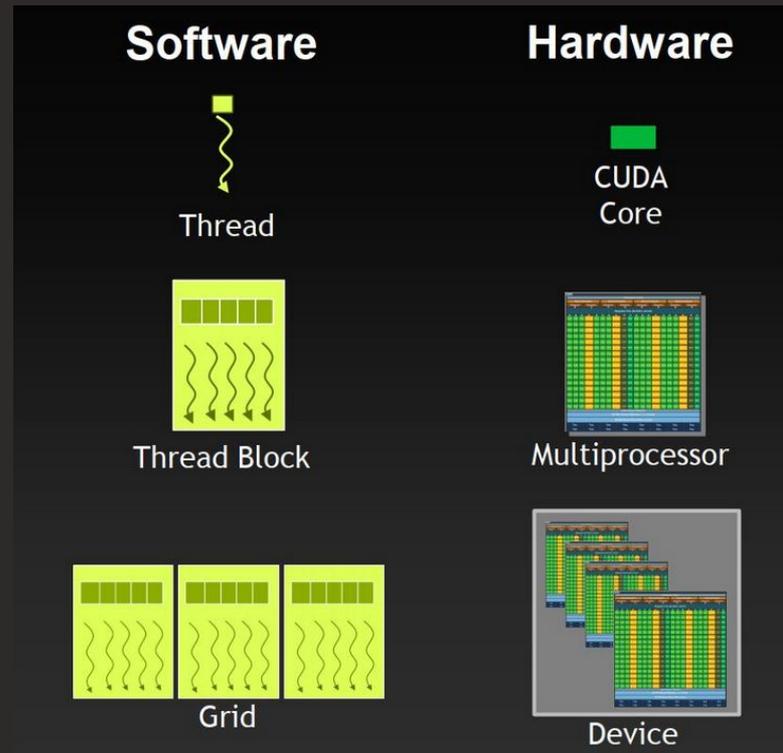
Modelo de programación GPU



Modelo de programación GPU #1

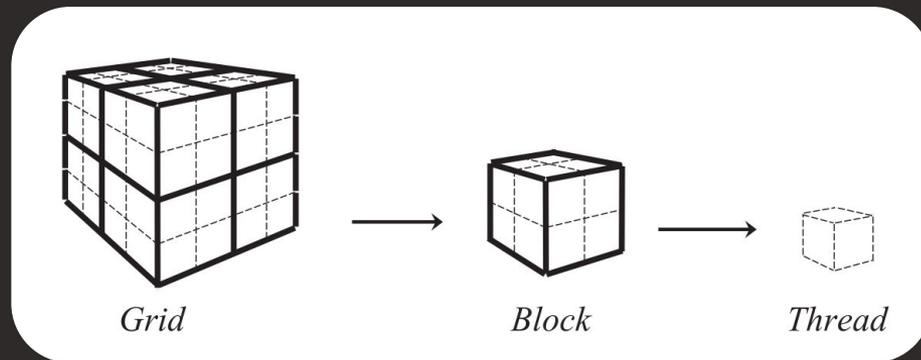


Modelo de programación GPU #2



Modelo de programación GPU #3

- El desarrollador se encarga de:
 - Mover datos Host \longleftrightarrow Device
 - Definir Grid, block y threads
 - 1D, 2D o 3D.
 - Escribir Kernel
 - Código *per-thread*.
 - *SIMD a nivel micro*.
 - Warp (32 threads)
 - Sincronos por bloque.
 - *SPMD a nivel macro*.
 - Blocks asíncronos



Modelo de programación GPU #4

Ejemplo CUDA: `main.cu`

```
__global__ void saxpy_kernel( float a, float* x, float* y, int n ){  
    int i;  
    i = blockIdx.x*blockDim.x + threadIdx.x;  
    if( i <= n ) x[i] = a*x[i] + y[i];  
}
```

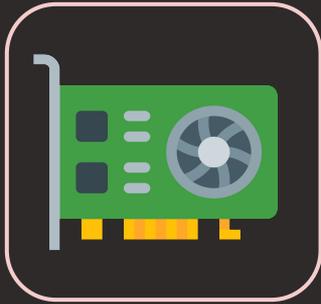
→ Parte Device (Kernel)

```
void saxpy( float a, float* x, float* y, int n ){  
    float *xd, *yd;  
    cudaMalloc( (void**)&xd, n*sizeof(float) );  
    cudaMalloc( (void**)&yd, n*sizeof(float) ); cudaMemcpy( xd, x,  
        n*sizeof(float),  
            cudaMemcpyHostToDevice );  
    cudaMemcpy( yd, y, n*sizeof(float),  
        cudaMemcpyHostToDevice );  
    saxpy_kernel<<< (n+31)/32, 32 >>>( a, xd, yd, n );  
    cudaMemcpy( x, xd, n*sizeof(float),  
        cudaMemcpyDeviceToHost );  
    cudaFree( xd ); cudaFree( yd );  
}
```

→ Parte Host

Compilacion:

```
nvcc -O3 main.cu -o prog
```



Introducción



Desafíos

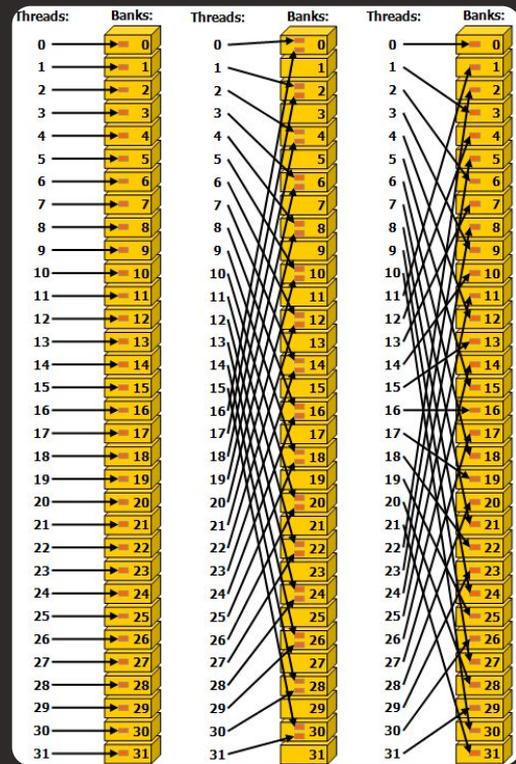
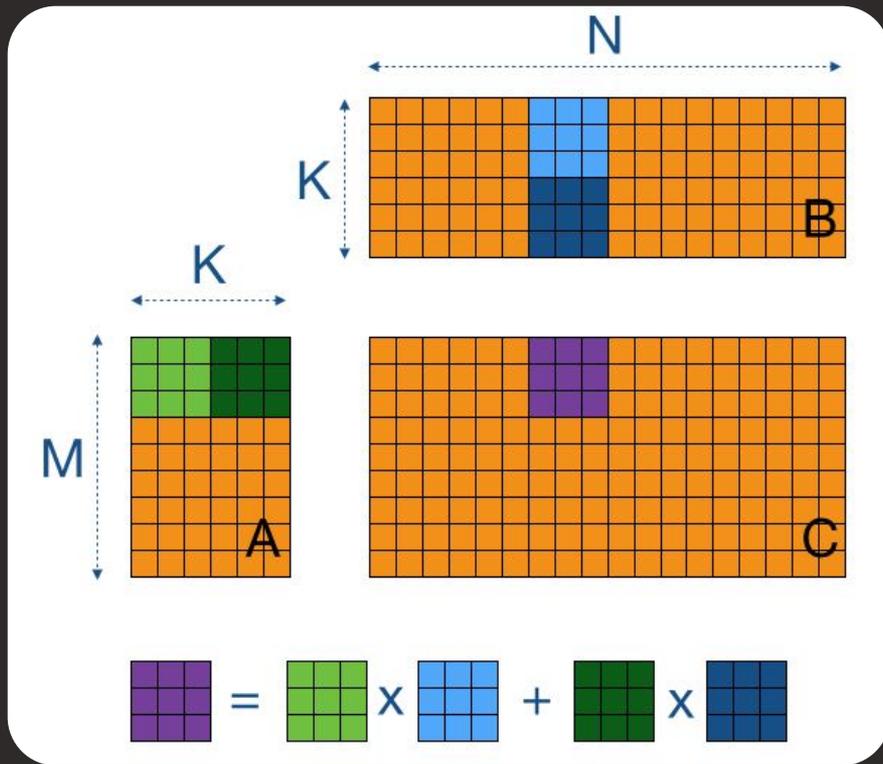


Grupo Temporal

Desafíos clásicos en GPU Computing

- **Desafíos de computación paralela:**
 - Encontrar casos de paralelización masiva
 - Algebra lineal (**matrix multiply**).
 - N-body problem.
 - Molecular dynamics.
 - PDEs (stencils).
 - Usos eficientes de caché L1 programmable.
 - Encontrar **localidad**
 - Evitar **bank conflicts**
 - Implementaciones **multi-GPU** eficientes.
 - Estructuras de datos jerárquicas eficientes en GPUs:
 - QuadTree, octree, BVH, kd-tree, k2-tree, etc
 - Usos exitosos de **dynamic-parallelism** (feature GPU)

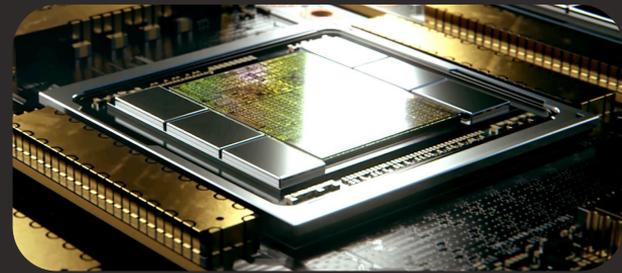
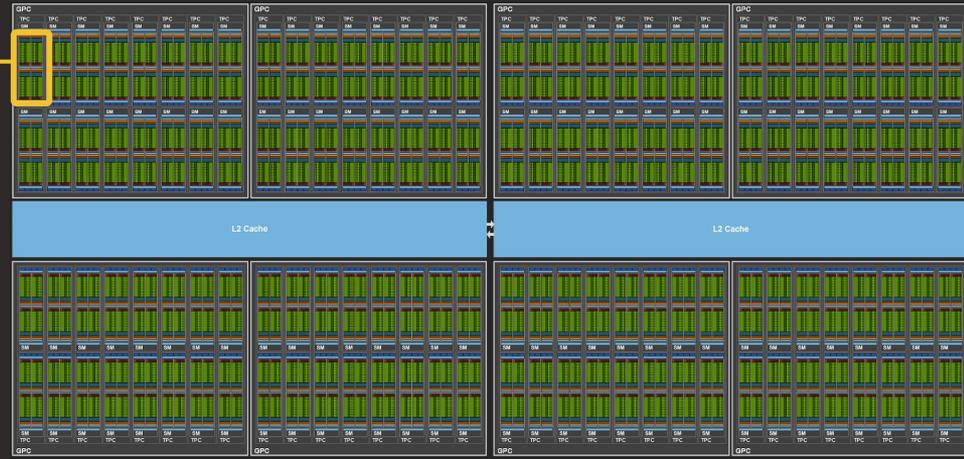
Usos eficientes de caché L1 programable



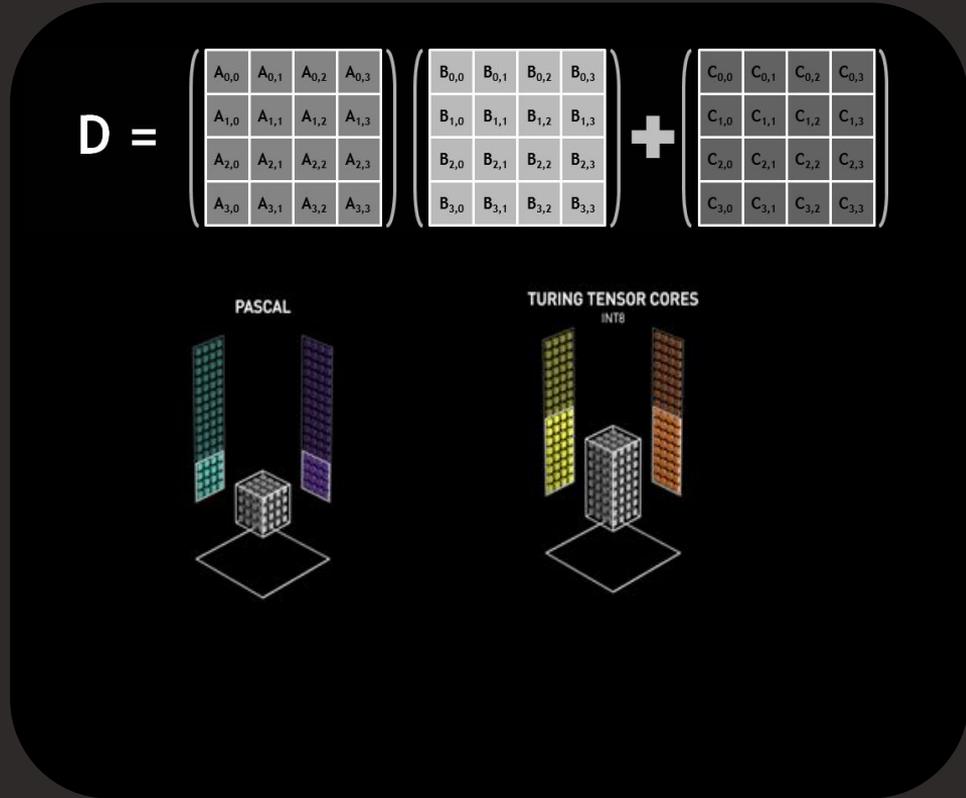
Desafíos del **presente** en GPU Computing

- Desarrollo de generadores / optimizadores de código paralelo
 - Full automaticos:
 - Semi-automaticos (directivas):
 - OpenACC, OpenMP, oneAPI
- Desarrollo y masificación de librerías / software:
 - cuBLAS
 - Thrust
 - Pytorch / TensorFlow
 - Gromacs
- Encontrar **nuevos usos a los ASICs** de GPU
 - Tensor-cores
 - RT cores.

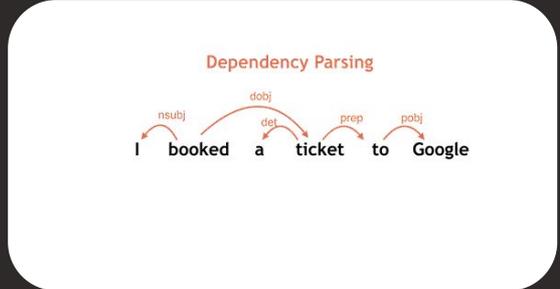
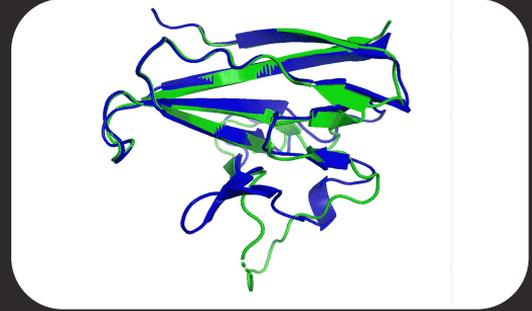
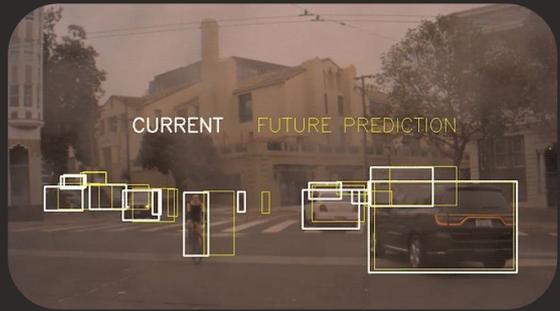
Tensor Cores



1 tensor core → 1 producto de matrices pequeñas (MMA)



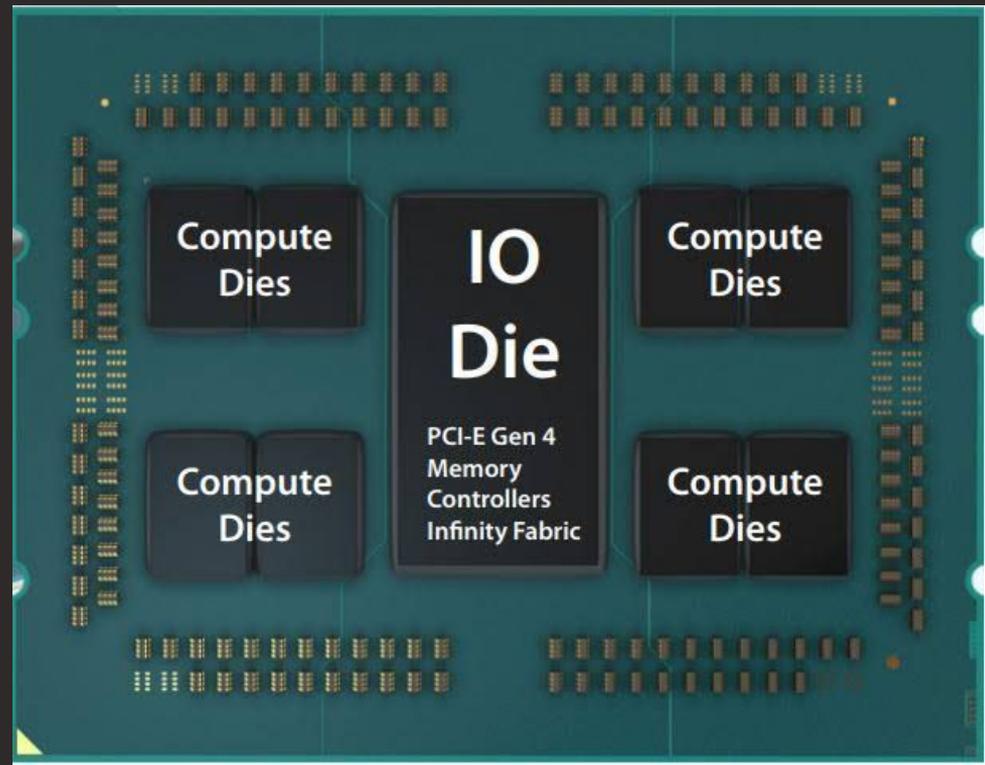
Uso natural de tensor cores

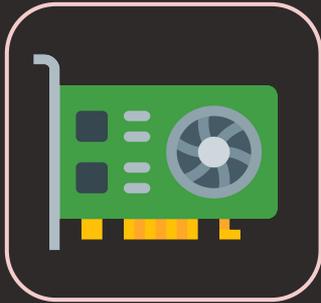


Desafíos del futuro en GPU Computing

- Desarrollar **algoritmos energéticamente eficientes**
 - Sacrificar un poco de rendimiento por eficiencia energética.
- **Combinar el uso de ASICs de GPU.**
 - Cores-regulares + tensor-cores + RT-cores.
 - Las aplicaciones guiarán el desarrollo de futuros ASICs.
 - Anticiparse a posibles nuevos ASICs.
- **Estudiar la eficiencia de los *multi-chip-module* (MCM)** en el diseño de nuevos algoritmos paralelos.

Multi-chip-module (MCM)





Introducción



Desafíos



Grupo Temporal

Grupo Temporal



Core Researchers



Dr. Cristóbal A. Navarro

Instituto de Informática, UACH, Chile



Dr. Nancy Hitchfeld

DCC, Universidad de Chile, Chile



Dr. Hector Ferrada

Instituto de Informática, UACH, Chile

Collaborators



Dr. Ricardo Barrientos

Catholic University of Maule, Chile



Dr. Miguel Romero

Adolfo Ibáñez University, Chile



Dr. Benoit Crespin

University of Limoges, France



Dr. Benjamin Bustos

DCC, Universidad de Chile, Chile



Dr. Mario Gonzalez

Instituto Acustica, UACH, Chile

Students



Roberto Carrasco

PhD Student at DCC, Uchile, Chile.



Felipe Quezada

MSc student at MIN, UACH, Chile.



Roberto Melita

MSc student at MIN, UACH, Chile.



Rodrigo Stevenson

MSc student at MIN, UACH, Chile.



Enzo Meneses

MSc student at MIN, UACH, Chile.



Carlos Schoenfeldt

Undergraduate, UACH, Chile.



Diego Rojas

Undergraduate, UACH, Chile.

<http://temporal.uach.cl>



- **Investigación en *GPU Computing*:**
 - Aceleración de algoritmos de geometría/simulación.
 - Mapeo eficiente de threads.
 - Paralelismo eficiente en memoria dinámica.
 - Eficiencia energética en GPUs.



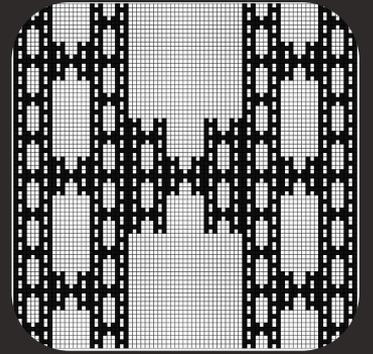
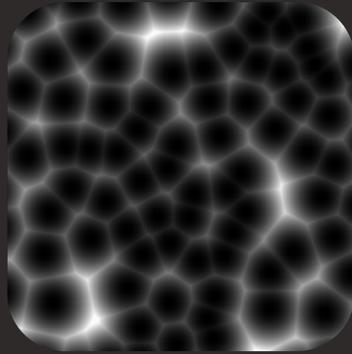
Universidad Austral de Chile
Conocimiento y Naturaleza



UNIVERSIDAD
DE CHILE

Puede el cómputo no-IA también beneficiarse de los tensor cores?

$$\sum_{i=1}^n x_i$$



[2018-2021] *Fondecyt iniciación N° 11180881*



Gracias

Repositorio Temporal: <https://github.com/temporal-hpc>

Investigación #1

Sumatoria por tensor core

$$\sum_{i=1}^n x_i$$

Idea clave #1: un par de MMAs pueden sumar m^2 datos.

$$D = J_m \times \begin{bmatrix} x_{11} & \dots & x_{1m} \\ \vdots & \ddots & \vdots \\ x_{m1} & \dots & x_{mm} \end{bmatrix} + 0_m \quad (9)$$

$$= \begin{bmatrix} \sum_{i=1}^m x_{1i} & \dots & \sum_{i=1}^m x_{1i} \\ \vdots & \ddots & \vdots \\ \sum_{i=1}^m x_{mi} & \dots & \sum_{i=1}^m x_{mi} \end{bmatrix} \quad (10)$$

MMA #1

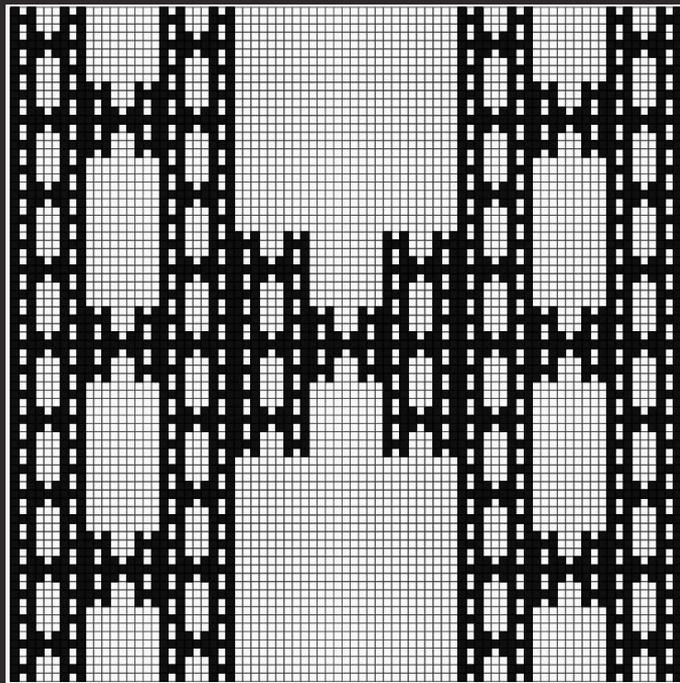
$$D' = \begin{bmatrix} \sum_{i=1}^m x_{1i} & \dots & \sum_{i=1}^m x_{1i} \\ \vdots & \ddots & \vdots \\ \sum_{i=1}^m x_{mi} & \dots & \sum_{i=1}^m x_{mi} \end{bmatrix} \times J_m + 0_m \quad (11)$$

$$= \begin{bmatrix} \sum_{i=1}^m \sum_{j=1}^m x_{ij} & \dots & \sum_{i=1}^m \sum_{j=1}^m x_{ij} \\ \vdots & \ddots & \vdots \\ \sum_{i=1}^m \sum_{j=1}^m x_{ij} & \dots & \sum_{i=1}^m \sum_{j=1}^m x_{ij} \end{bmatrix} \quad (12)$$

MMA #2

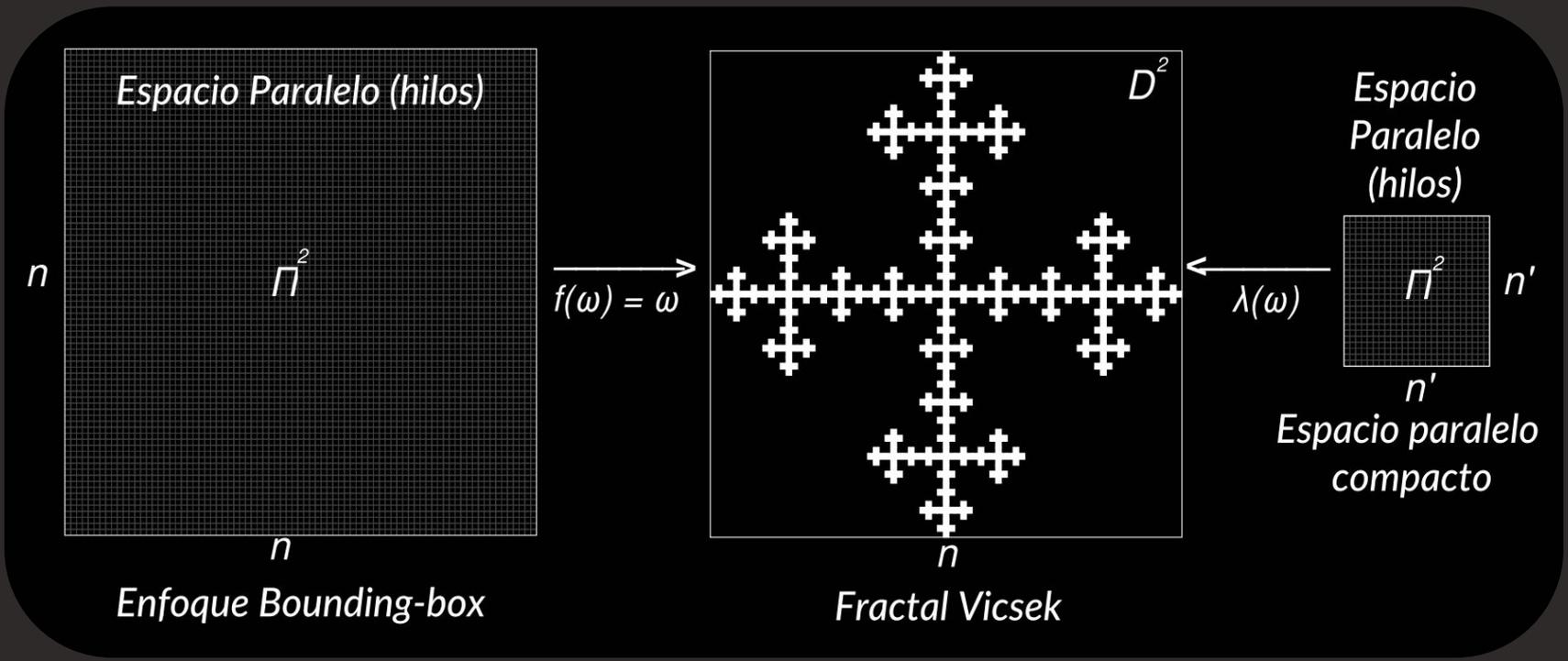
Investigación #2

Mapeo de *threads* de GPU
a geometría fractal discreta



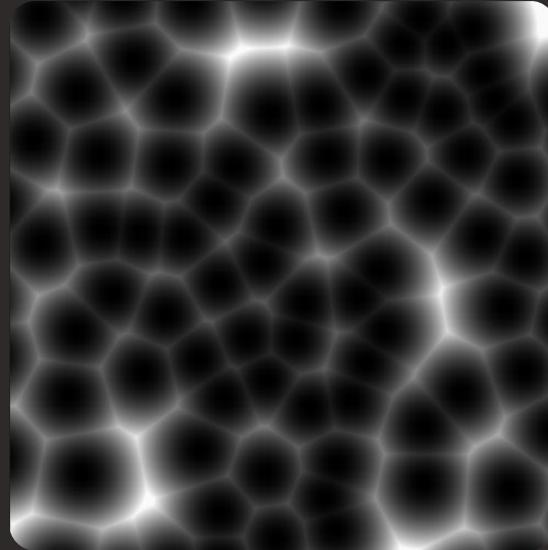
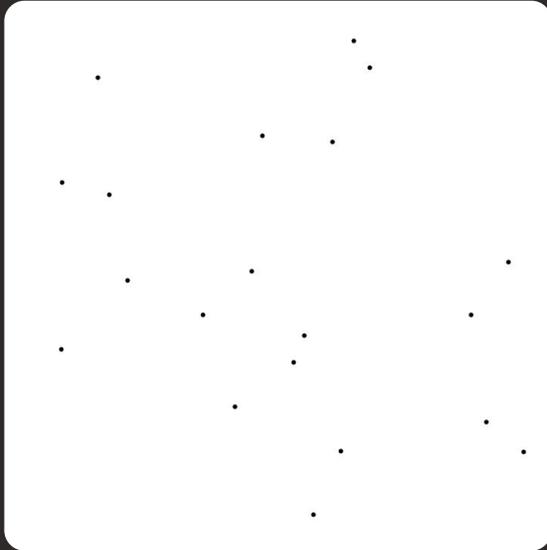
Idea estandar

Idea propuesta



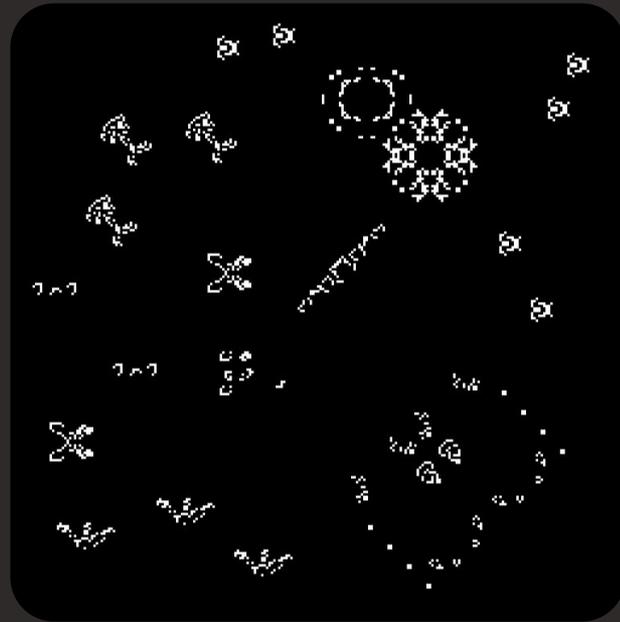
Investigación #3 (en curso)

Diagramas de Voronoi dinámicos acelerados por GPU



Investigación #4 (en curso)

Simulación de autómatas celulares con tensor cores



Investigación #6 (en curso)

Caracterización de la eficiencia energética de GPUs.

- Resultados preliminares confirman alta eficiencia energética con GPU.
- Actualmente obteniendo resultados experimentales con/sin tensor cores.

