# Judging a Type by Its Pointer: Optimizing GPU Virtual Functions
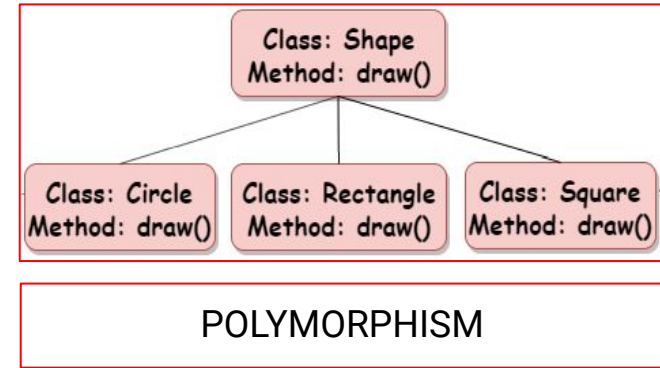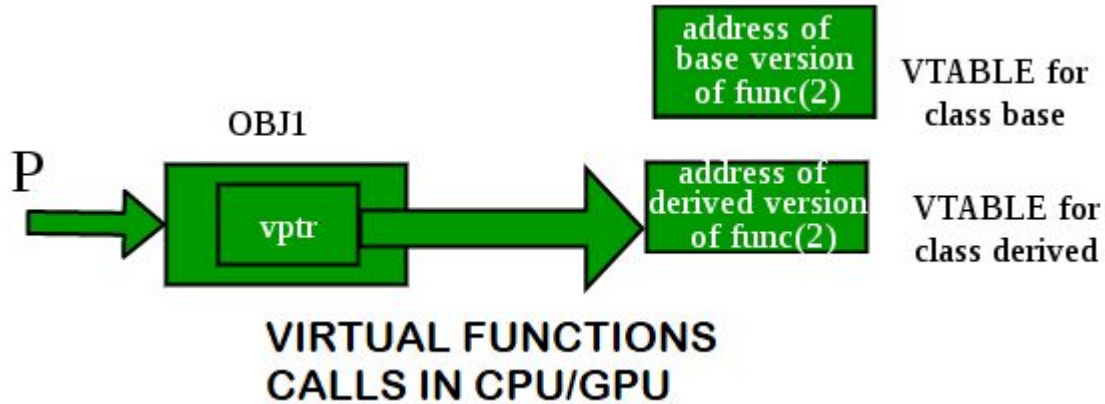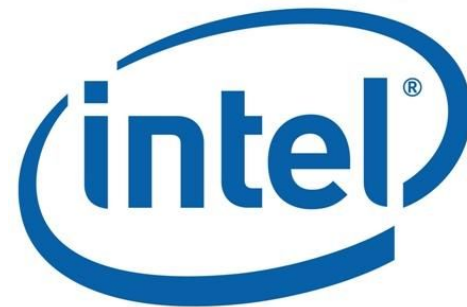
Patrick Xavier Márquez Choque
Jean Carlo Cornejo Cornejo

# Introducción



address of
base version
of func(2)

VTABLE for
class base

OBJ1

P

vptr

address of
derived version
of func(2)

VTABLE for
class derived

**VIRTUAL FUNCTIONS
CALLS IN CPU/GPU**

Class: Shape
Method: draw()

Class: Circle
Method: draw()

Class: Rectangle
Method: draw()

Class: Square
Method: draw()

POLYMORPHISM

# Problema del Paper









No Research in GPU's Applications

SharedOA
Framework

# Propuestas del Paper



COAL

&

TypePointer

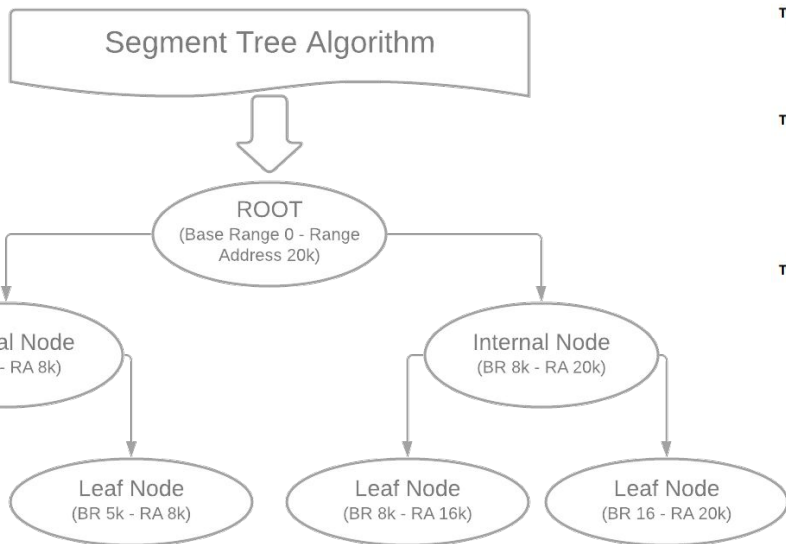Coordinated Object Allocation
and function Lookup

Hardware Support MMU Mod

# COAL



Segment Tree Algorithm

ROOT
(Base Range 0 - Range Address 20k)

Internal Node
(BR 0 - RA 8k)

Internal Node
(BR 8k - RA 20k)

Leaf Node
(BR 0 - RA 5k)

Leaf Node
(BR 5k - RA 8k)

Leaf Node
(BR 8k - RA 16k)

Leaf Node
(BR 16 - RA 20k)

**Algorithm 1:** Scan algorithm for the virtual range table

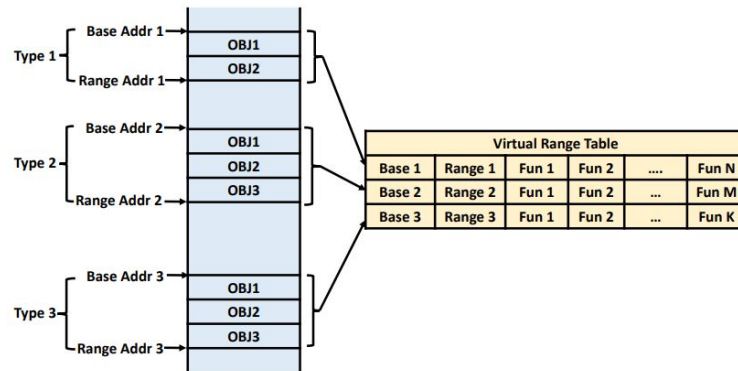**Function** ObjectRangeLookup($objectAddr$, $funcIndex$)
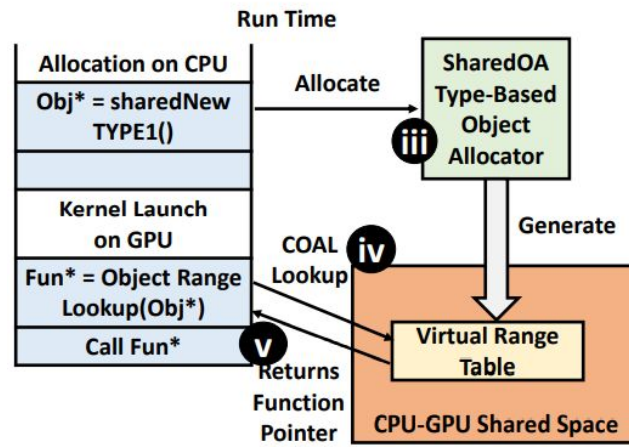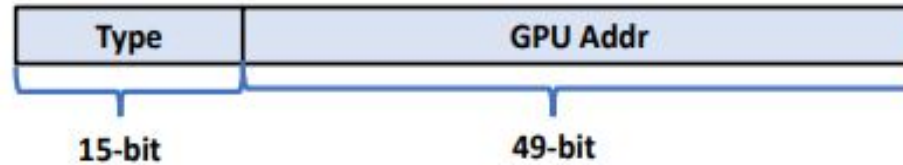
## Complexity: O(log$_2$ (K))



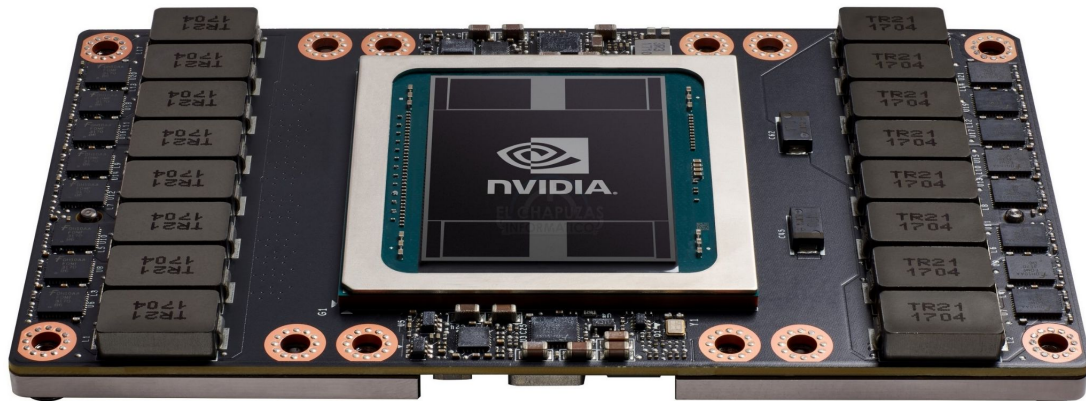Figure 3: Type-based object allocator.

# TypePointer

**TypePointer Format**

| Type | GPU Addr |
|------|----------|
| 15-bit | 49-bit |

1: SHR     $Ra,$    $R_{obj},$    #49

2: ADD     $Ra,$    $Ra,$    $R_{vTablesStartAddr}$

3: LDG     $Ra,$    $\#vFuncOffset\ (Ra)$

4: CALL     $Ra$

**(b) Instructions implementing virtual function calls with *Type-Pointer***

# Metodología



TypePointer

GPGPU-SIM

Hardware de Pruebas, con limitación de 4GB como HEAP size máximo

Workloads: Dynasoar, GraphChi-vE, Open Source Ray Tracer

CUDA   Concord   SharedOA   COAL   TypePointer

# Experimentación

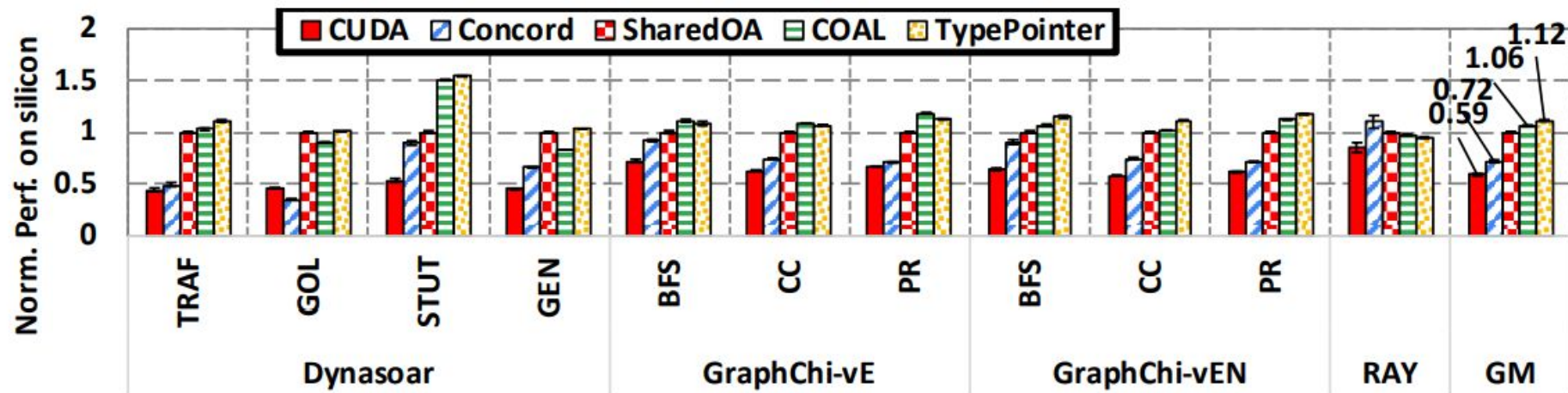| Workload | Description | # Objects | # Types | vFuncs |
|---|---|---|---|---|
| **Dynasoar Workloads [22, 46]** | | | | |
| Traffic (TRAF) | A Nagel-Schreckenberg model traffic simulation to model streets, cars and traffic lights for traffic flows. | 1573714 | 6 | 74 |
| Game Of Life (GOL) | Game of life is a cellular automaton formulated by John Horton Conway. This benchmark has two abstract classes Cells and Agent. | 5645916 | 4 | 29 |
| Structure (STUT) | Structure uses the Finite element method to simulate the fracture in a material. The benchmark models the material with springs and nodes. | 525000 | 4 | 40 |
| Generation (GEN) | Generation is an extension of gol benchmark. The cells in Generation have more intermediate states which lead to more complicated scenarios. | 1048576 | 4 | 33 |
| **GraphChi-vE Workloads [35]** | | | | |
| Breadth First Search (BFS) | BFS traverses graph nodes and updates a level field in a breadth-first manner. The GraphChi-vE BFS implementation defines an abstract class for edges, ChiEdge, and a concrete classEdge, which implements all the virtual functions of ChiEdge. | 2254419 | 4 | 5 |
| Connected Components (CC) | Connected Component is commonly used for image segmentation and cluster analysis, it employs an iterative node updates according to the labels of adjacent nodes. | 2254419 | 4 | 6 |
| Page Rank (PR) | Page rank is a classic algorithm to rank the pages of search engine results using iterative updates for each node. | 2254419 | 4 | 3 |
| **GraphChi-vEN Workloads [36]** | | | | |
| Breadth First Search (BFS) | The GraphChi-vEN BFS implementation also defines an abstract base class for vertex, ChiVertex, and a concrete class vertex, which implements ChiVertex's virtual functions. | 2254419 | 4 | 15 |
| Connected Components (CC) | GraphChi-vEN CC is similar to GraphChi-vE described above. However, GraphChi-vEN CC has both virtual edges and nodes. | 2254419 | 4 | 15 |
| Page Rank (PR) | GraphChi-vEN PR is similar to GraphChi-vE described above. However, GraphChi-vEN PR has both virtual edges and nodes. | 2254419 | 4 | 10 |
| **Open Source Ray Tracer [40]** | | | | |
| Raytracing (RAY) | RAY performs global rendering of of spheres and planes. The algorithm traces light rays through a scene, bouncing them off of objects, and back to the screen. | 1000 | 3 | 3 |

# Resultados



Figure 6: Performance, normalized to *SharedOA* on a silicon V100 GPU, averaged over 10 runs (error-bars=max and min).
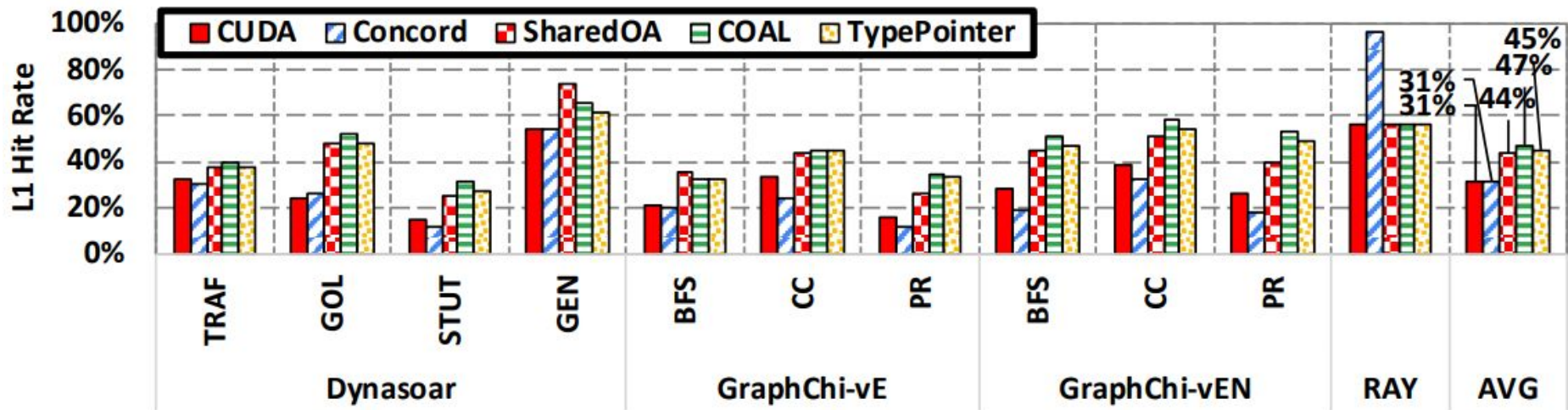
# Resultados



Figure 9: L1 cache hit rate on a silicon V100 GPU.

# Conclusiones



PERFORMANCE:
80%, 47% & 6% IMPROVEMENT



PERFORMANCE:
90%, 56% & 12%
IMPROVEMENT

GREAT WORK!

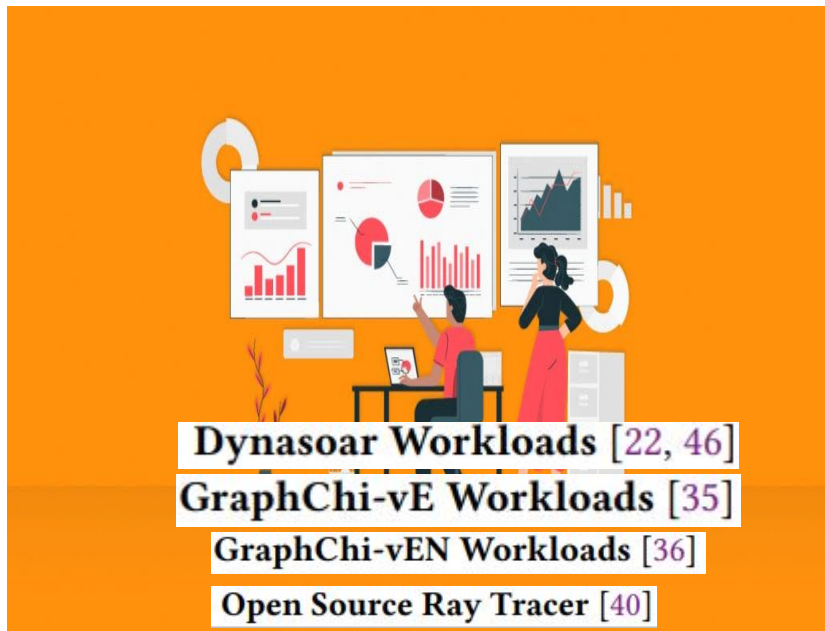🟥 CUDA   🟦 Concord   🟥 SharedOA   🟩 COAL   🟨 TypePointer

# Descripción de las dependencias para las réplicas del trabajo

- **Hardware Dependencies**: Intel x86 machine with NVIDIA Volta architecture GPU with at least 16GB GPU memory. The experiments in this paper use V100 GPU with 32GB GPU memory. 1GB CPU memory is required to contain the evaluation repositories.

- **Software Dependencies**. Ubuntu 18.04 Linux is preferred to run the experiments. NVIDIA CUDA 10.1, 10.2 and 11.1 and GPU driver are required to compile and run the GPU workloads.

# Experimentación

Utilización de los Workflow( Benchmarking ):

Experimentos con ejemplos simples:

Dynasoar Workloads [22, 46]
GraphChi-vE Workloads [35]
GraphChi-vEN Workloads [36]
Open Source Ray Tracer [40]

## SharedOA, COAL and TypePointer

### Software prerequisite

- Ubuntu 18.04.5 LTS Linux
- git
- Python 2.7
- CUDA 10.1: https://developer.nvidia.com/cuda-10.1-download-archive-base
- Transform script for COAL and TypePointer: https://github.com/brad-mengchi/asplos_2021_ae

# Experimentación con Ejemplos Simples (COAL)

```cpp
3   #include "COAL.h"
4   #define NUM_OBJ 512
5   class S1 {
6   public:
7     int var;
8     __host__ __device__ S1() {}
9     virtual __host__ __device__ void inc() = 0;
10    virtual __host__ __device__ void dec() = 0;
11  };
12
13  class S2 : public S1 {
14  public:
15    __host__ __device__ S2() {}
16    __host__ __device__ void inc() { this->var += 2; }
17
18    __host__ __device__ void dec() { this->var -= 2; }
19  };
```

```cpp
21  __global__ void kernel(S1 **ptr) {
22    int tid = threadIdx.x + blockDim.x * blockIdx.x;
23    // this variable must be defined in every kerenl that uses COAL
24    void **vtable;
25    if (tid < NUM_OBJ) {
26      COAL_S1_inc(ptr[tid]);
27      ptr[tid]->inc();
28    }
29  }
```

```cpp
50      my_obj_alloc.toDevice();
51      int blockSize = 256;
52      int numBlocks = (NUM_OBJ + blockSize - 1) / blockSize;
53      kernel<<<numBlocks, blockSize>>>(ptr);
54      cudaDeviceSynchronize();
```

# Réplica del Experimento en COAL

Réplica con el Cluster de Google Colab:

[Google Colab Experiment](#)

# Experimentación con Ejemplos Simples (SharedOA)

```cpp
#include "mem_alloc.h"
#define NUM_OBJ 512
class S1 {
public:
  int var;
  __host__ __device__ S1() {}
  virtual __host__ __device__ void inc() = 0;
  virtual __host__ __device__ void dec() = 0;
};

class S2 : public S1 {
public:
  __host__ __device__ S2() {}
  __host__ __device__ void inc() { this->var += 2; }

  __host__ __device__ void dec() { this->var -= 2; }
};
```

```cpp
__global__ void kernel(S1 **ptr) {
  int tid = threadIdx.x + blockDim.x * blockIdx.x;
  if (tid < NUM_OBJ)
    ptr[tid]->inc();
}
```

```cpp
// to access the virtual function from device we call
my_obj_alloc.toDevice();
int blockSize = 256;
int numBlocks = (NUM_OBJ + blockSize - 1) / blockSize;
kernel<<<numBlocks, blockSize>>>(ptr);
cudaDeviceSynchronize();
printf("Device Call Done\n");
```

# Réplica del Experimento en SharedOA

Réplica en Google Colab:



```
vtable [2S2][0]:0x8
vtable [2S2][1]:0x10
Objects Creation Done
Host Call Done
Device Call Done
Host Call Done
Device Call Done
ptr[0].var = 0
ptr[1].var = 0
ptr[2].var = 0
ptr[3].var = 0
ptr[4].var = 0
ptr[5].var = 0
ptr[6].var = 0
ptr[7].var = 0
ptr[8].var = 0
ptr[9].var = 0
```

```
ptr[508].var = 0
ptr[509].var = 0
ptr[510].var = 0
ptr[511].var = 0
```

```
Type#0:
Type Size: 16     Number of Buckets : 1    Range Size: 1048576    Number of Objs : 512

Avg Types Size 16.000000
```

# Experimentación con Ejemplos Simples (TP)

```
#include "mem_alloc_tp.h"
#include "TP.h"
#define NUM_OBJ 512
class S1 {
public:
  int var;
  __host__ __device__ S1() {}
  virtual __host__ __device__ void inc() = 0;
  virtual __host__ __device__ void dec() = 0;
};

class S2 : public S1 {
public:
  __host__ __device__ S2() {}
  __host__ __device__ void inc() { this->var += 2; }

  __host__ __device__ void dec() { this->var -= 2; }
};
```

```
__global__ void kernel(S1 **ptr) {
  int tid = threadIdx.x + blockDim.x * blockIdx.x;
  // this variable must be defined in every kerenl that uses COAL
  void **vtable;
  if (tid < NUM_OBJ) {
    S1 * obPtr = ptr[tid];
    TP_S1_inc(obPtr);
    CLEANPTR(obPtr,S1 *)->inc();
  }
}
```

```
my_obj_alloc.toDevice();
int blockSize = 256;
int numBlocks = (NUM_OBJ + blockSize - 1) / blockSize;
kernel<<<numBlocks, blockSize>>>(ptr);
cudaDeviceSynchronize();
```

# Réplica del Experimento en TP

Réplica en Google Colab:

```
vtable [2S2][0]:0x8
vtable [2S2][1]:0x10
total_count 1048576 , total 1048576 , type_size 16
Objects Creation Done
Type0:  (nil)
Type1:  0x8
#############
gpuptr 0x7f78dc001138
Device Call Done
ptr[0].var = 2
ptr[1].var = 2
```

```
ptr[510].var = 2
ptr[511].var = 2
Type#0:
Type Size: 16     Number of Buckets : 1    Range Size: 1048576     Number of Objs : 512

Avg Types Size 16.000000
```

# Link del Repositorio

[Github](#)

[Paper](#)