

# **Towards**

# **Autonomous**

# **Compiler Design Using Machine Learning**

**Zheng Wang**

School of Computing  
University of Leeds

*Work in collaboration with*

Chris Cummins

Zhanyong Tang

Pavlos Petoumenos

Hugh Leather

Guixin Ye

Huanting Wang

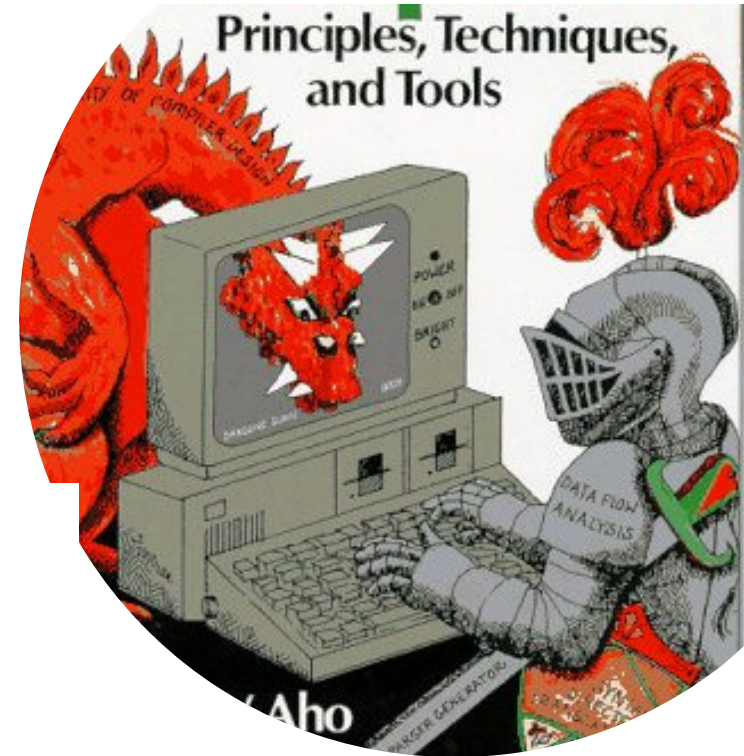
William Ogilvie

# Better compiler = happier users

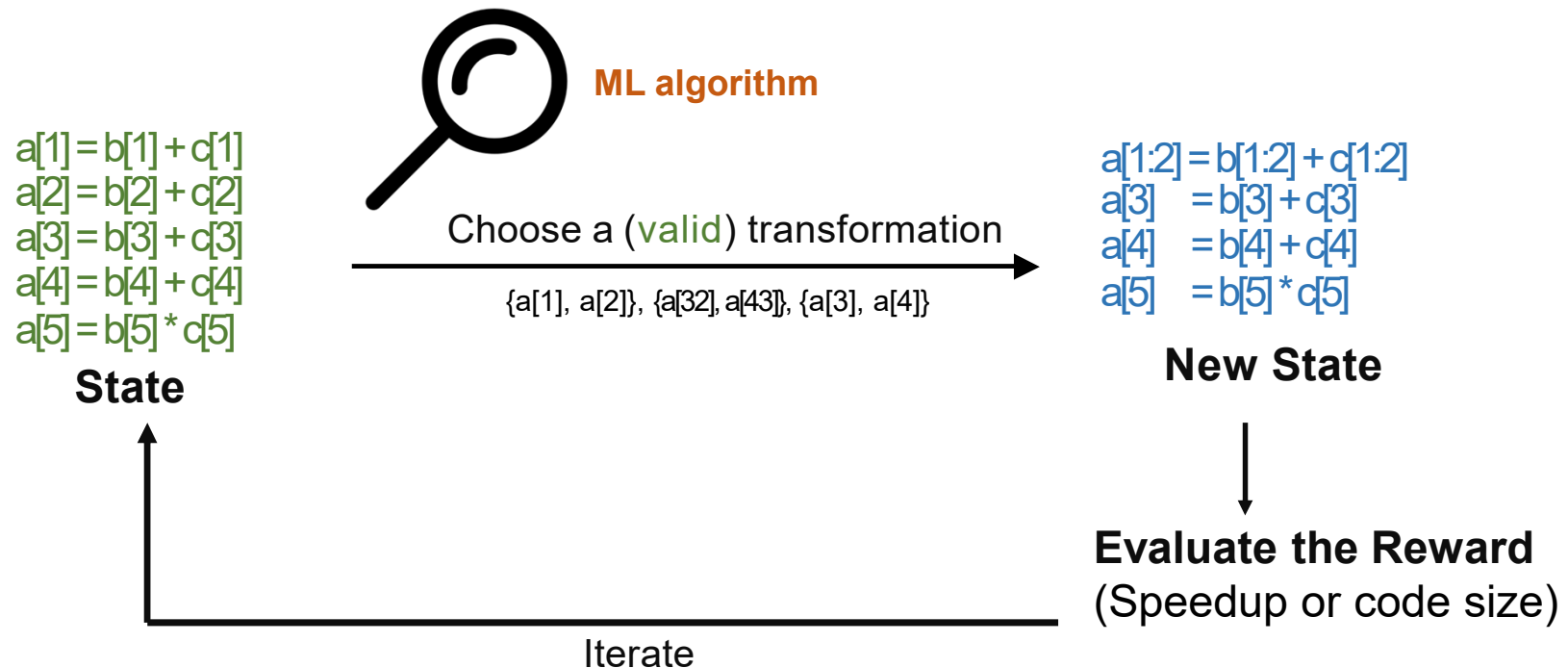
- Faster binary = better user experience 😊
- Lower hardware requirement or power = 💷 saved

Writing compiler optimisation heuristics **by hand** is time-consuming and hard

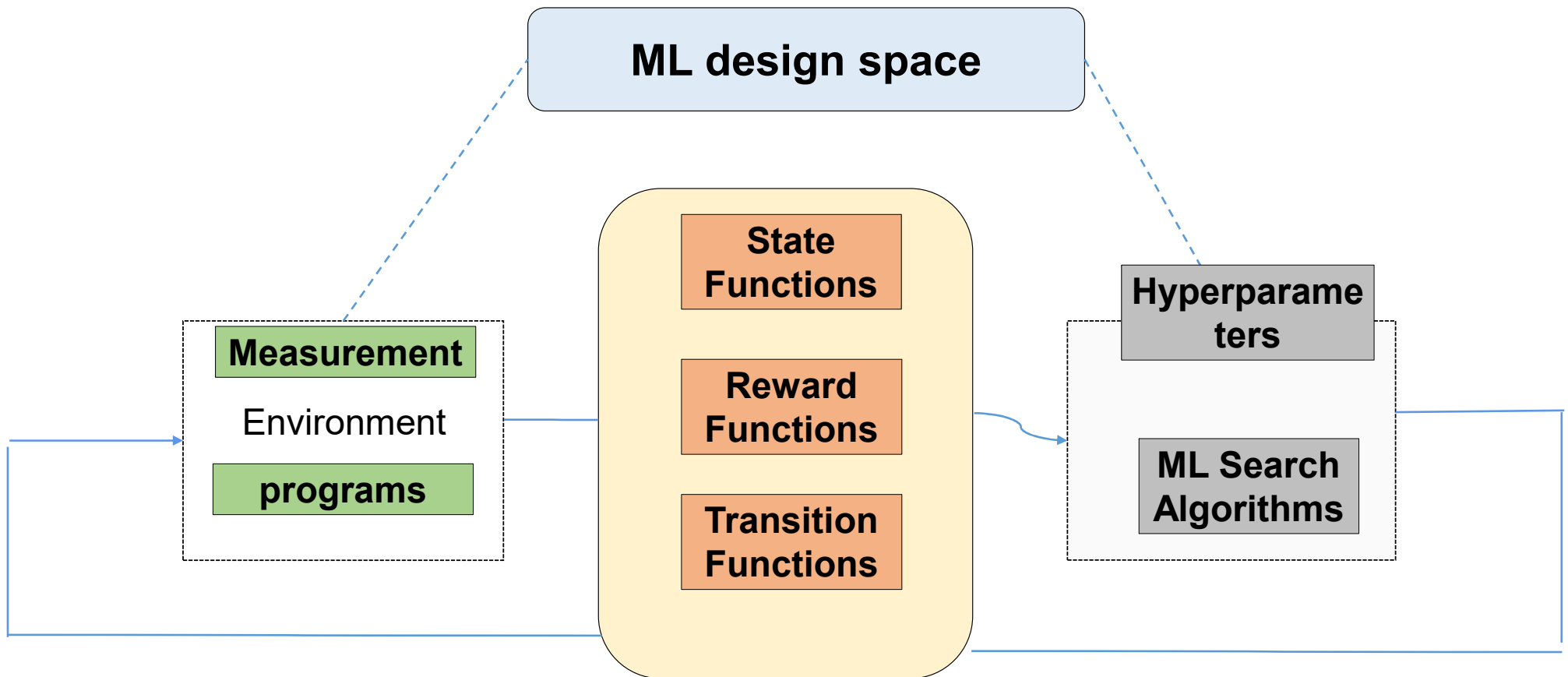
- **Machine learning** promises:  
    **Better compilers**  
    **Less development cost**



# Learn to search the optimisation space

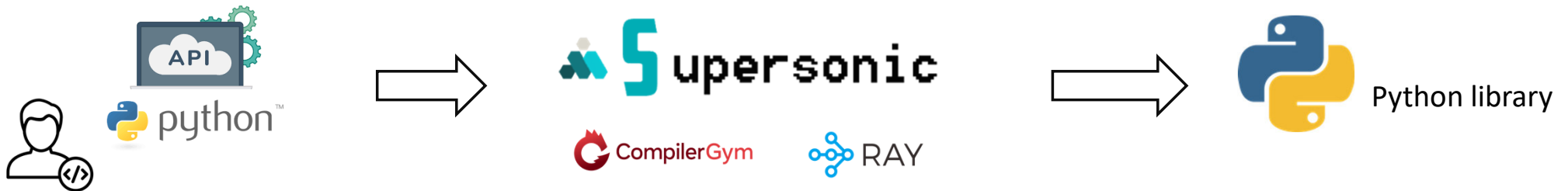


# ML components for compiler search



# Supersonic: AutoML for compiler optimisation

Lowering the barrier of integrating ML into compilers



Developers describes the compiler problem

Automatically selected and Tuned ML components

# User defines the search space

```
import Supersonic as ss
```

```
actions=["-O3", ... ]
```

Transformation options

```
class task(ss.PolicyInt):
```

```
    def __init__(self, benchmarks, *arg):  
        #Initialise an environment
```

Benchmarks for tuning

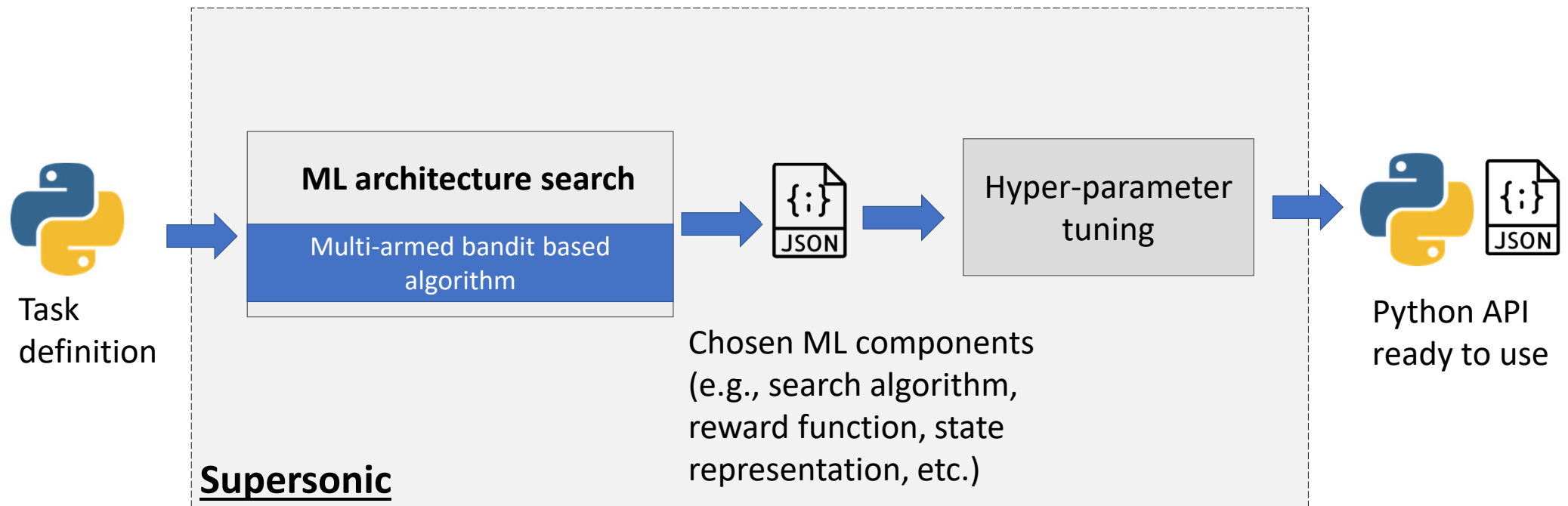
```
    def run(self):  
        #How to execute the compiled code
```

Measurement interface

```
    def step(self, code, action):  
        #take an action to transform/compile the code
```

How to apply an action

# Automatically find and tune the ML architecture





# Case studies

## Optimizing Image Pipelines (Halide)

➤ 4x prior methods

## Neural Network Code Generation (TVM)

➤ 7x prior methods

## Code Size Reduction

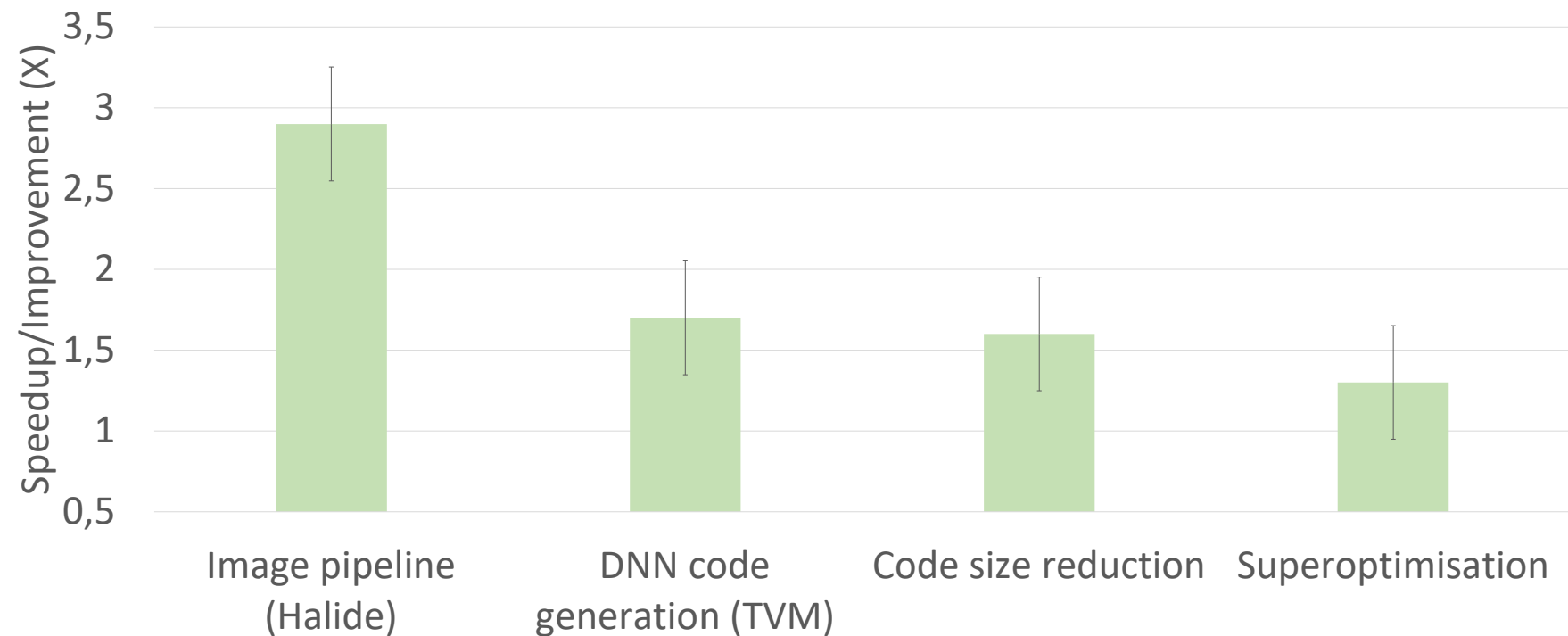
➤ 4 prior methods

## Superoptimization

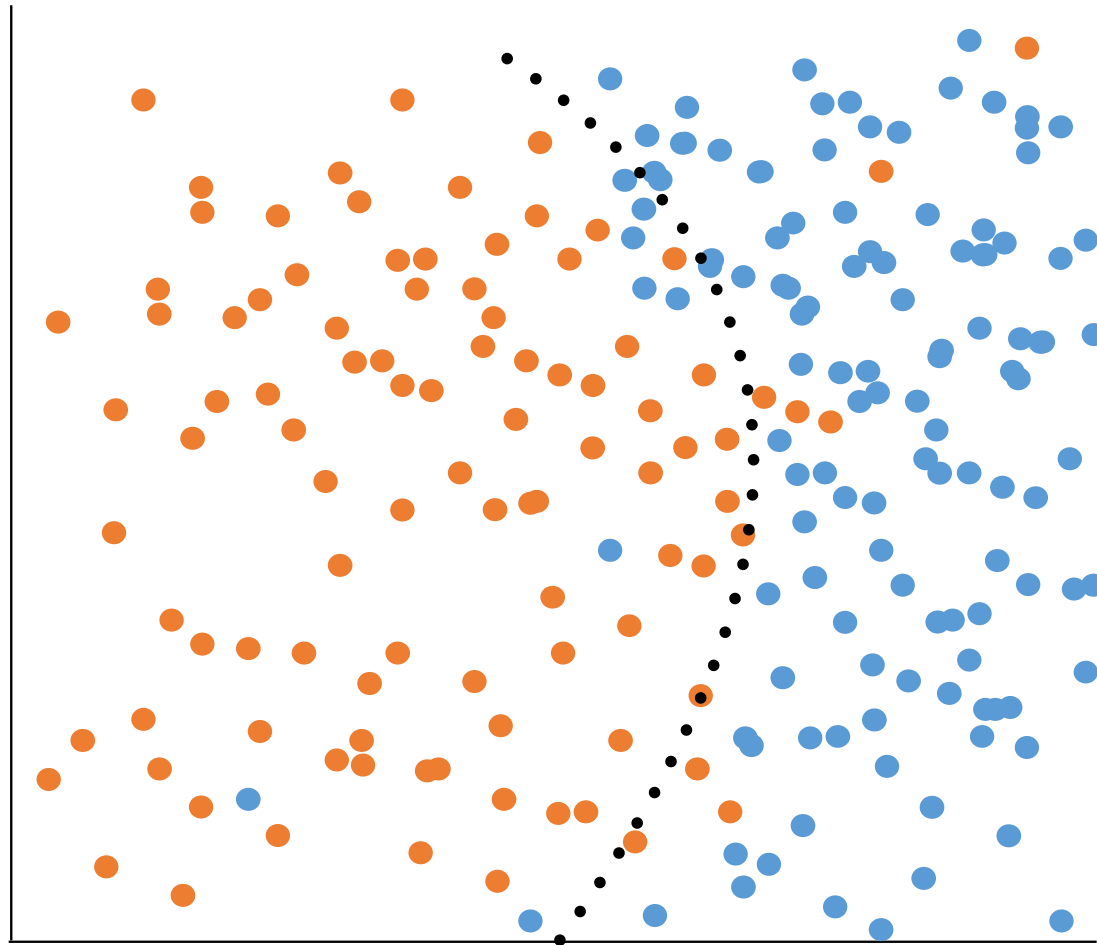
➤ 3 prior methods

Compare to hand-tuned ML approaches

# Even better than hand-crafted ML



# ML for predictive modelling



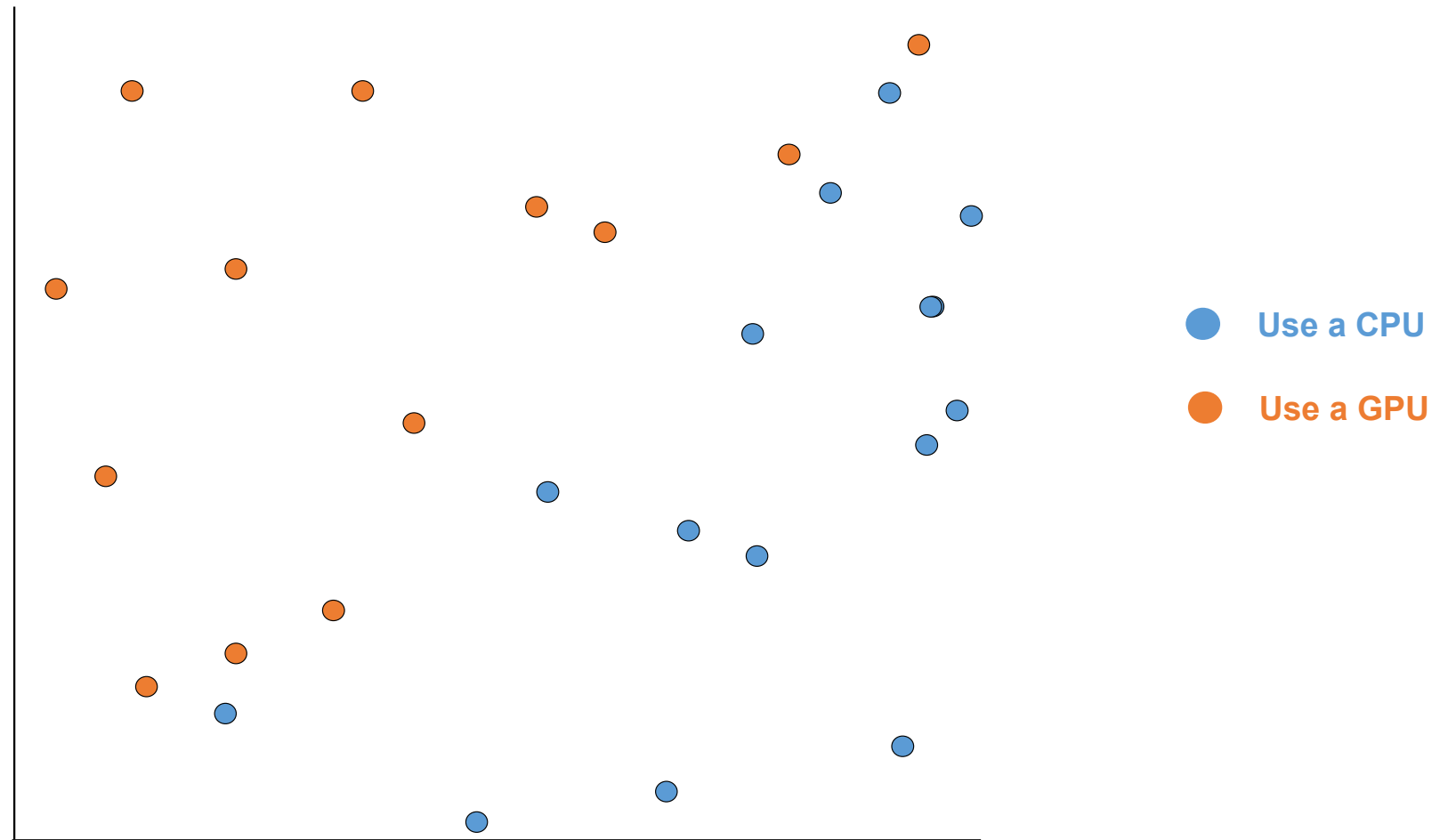
Which  
processor is  
faster?

● Use a CPU

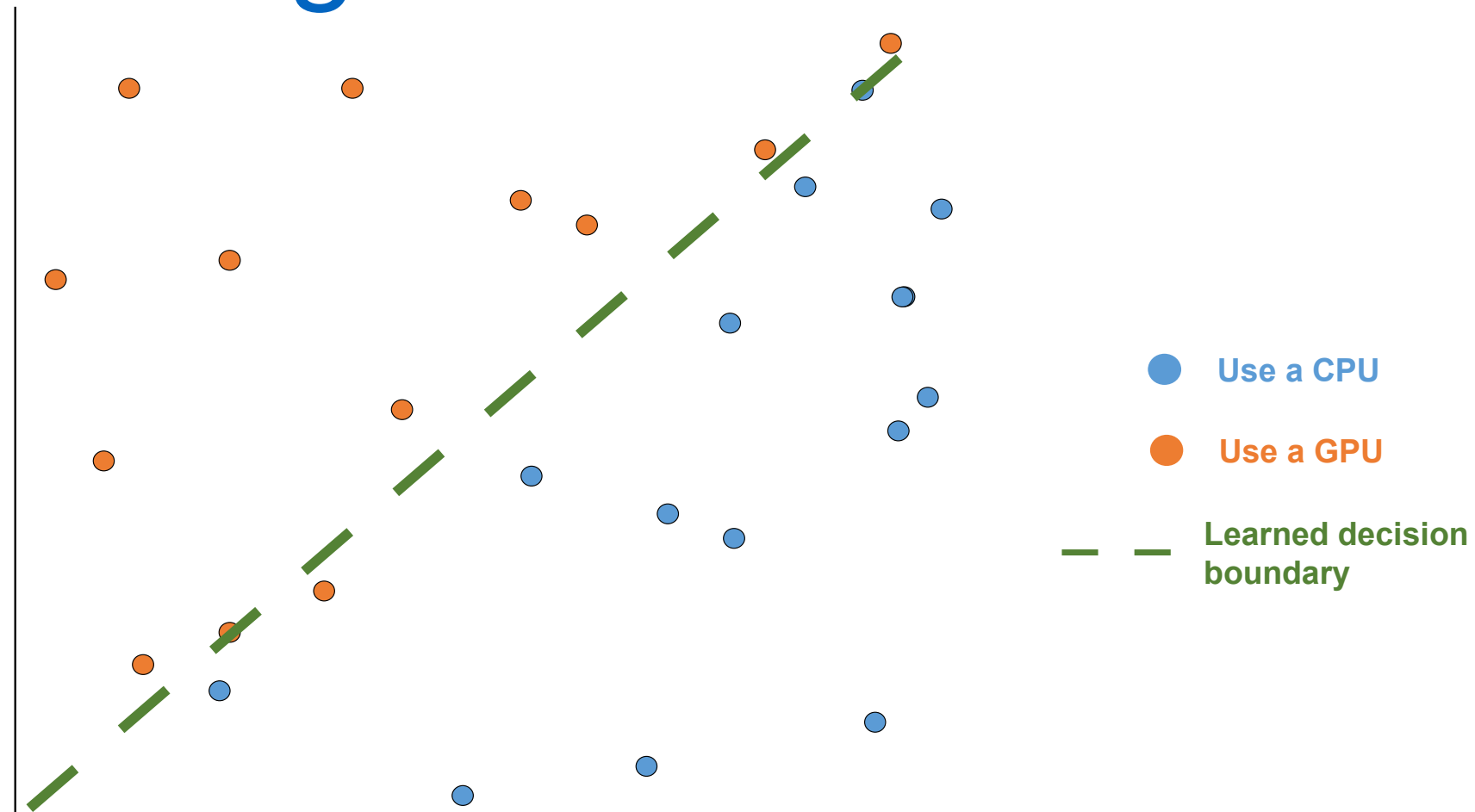
● Use a GPU

.... Target decision  
boundary

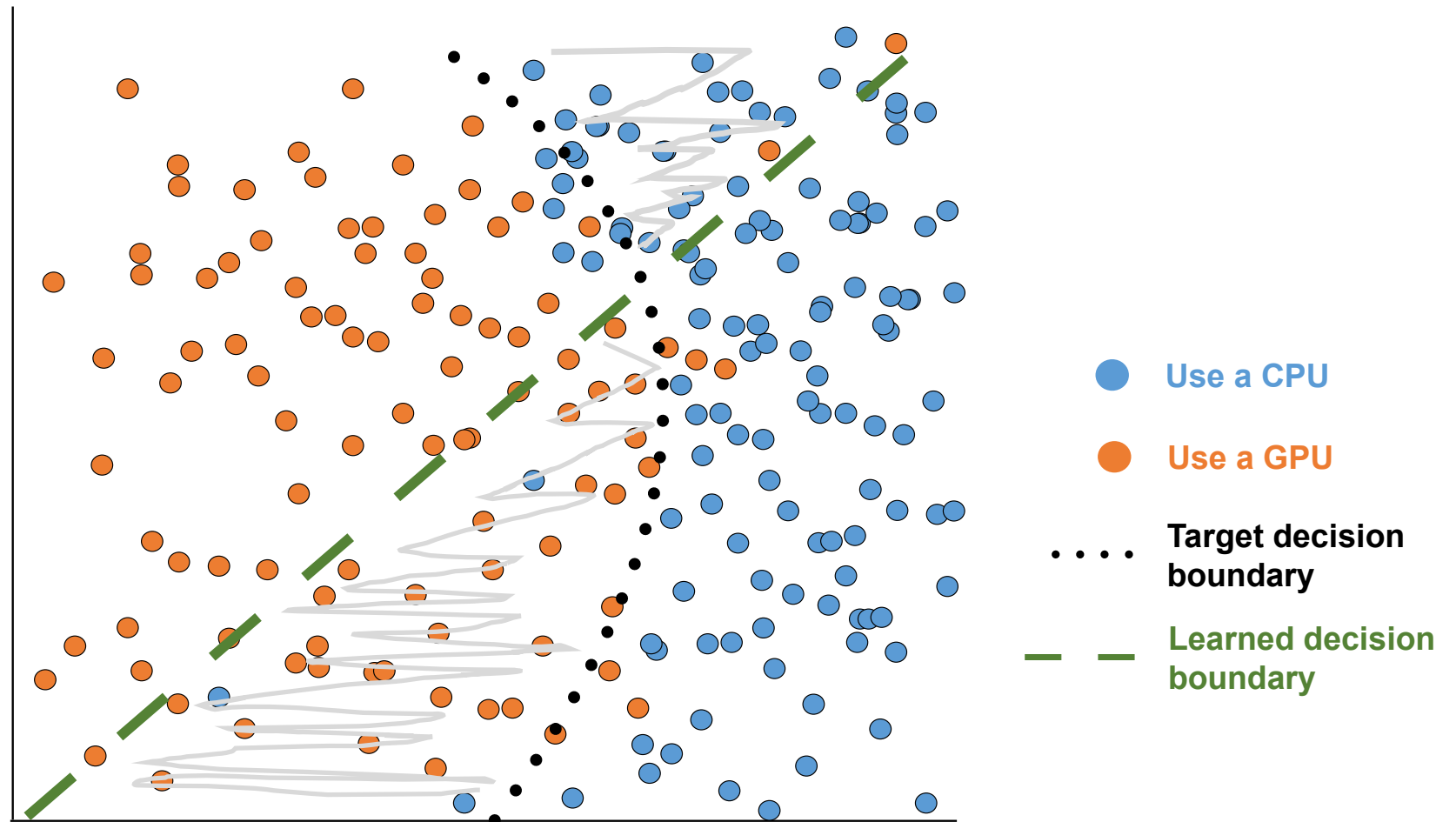
# We actually have:



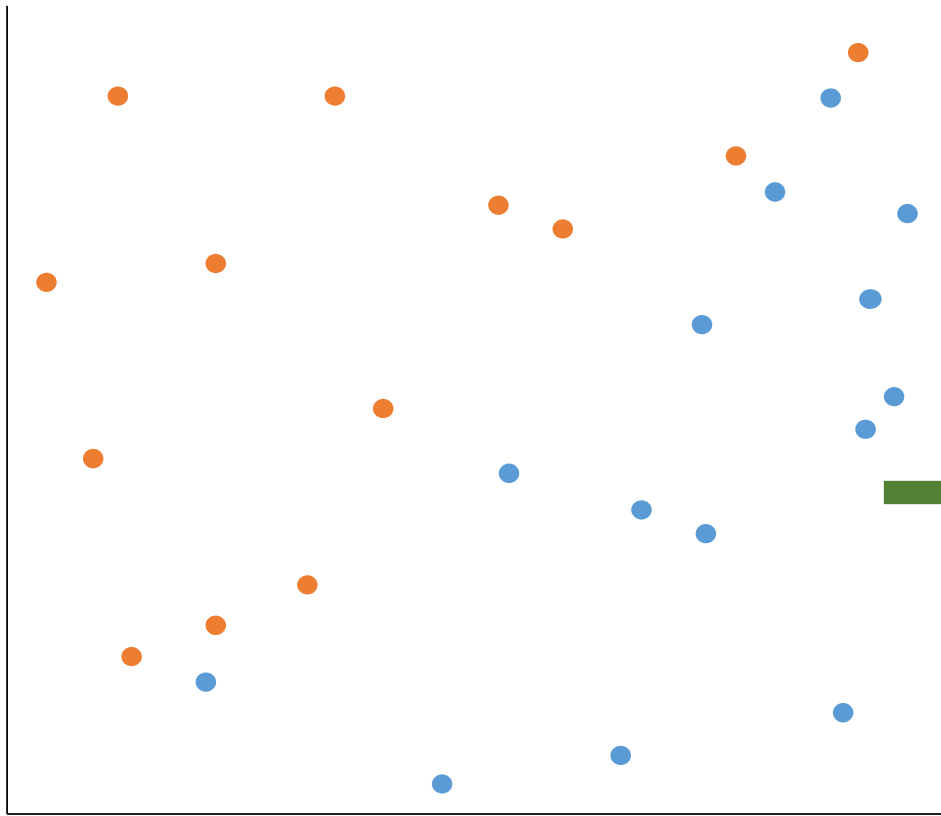
# This would give us a model like:



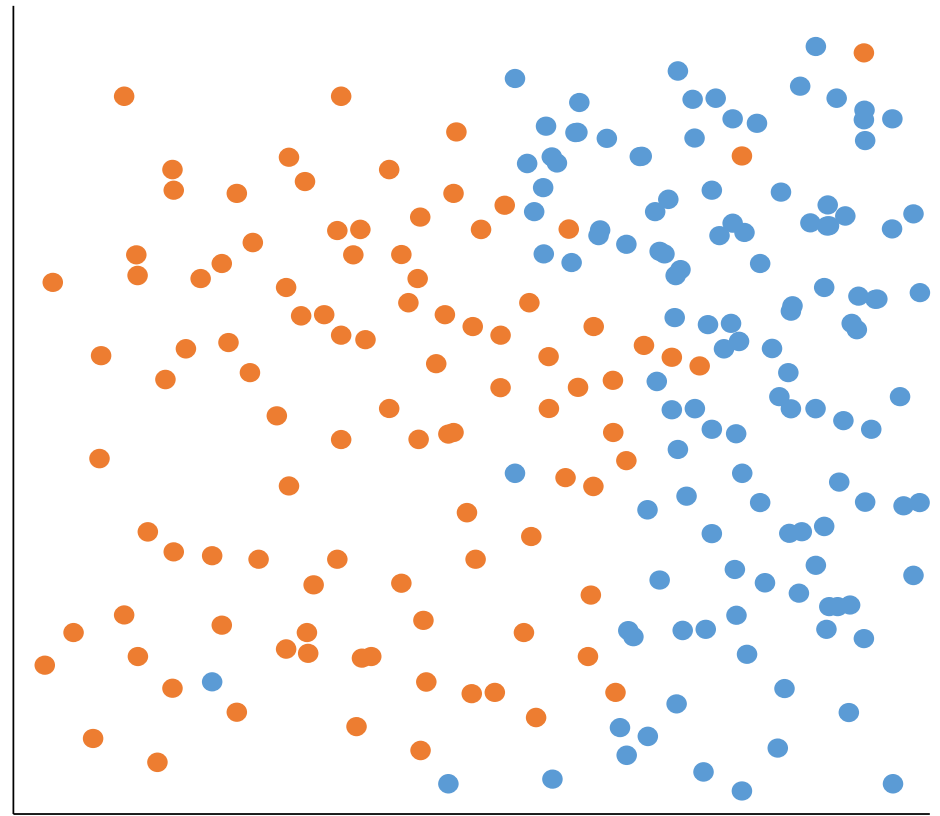
# *Insufficient training data give an inaccurate model*



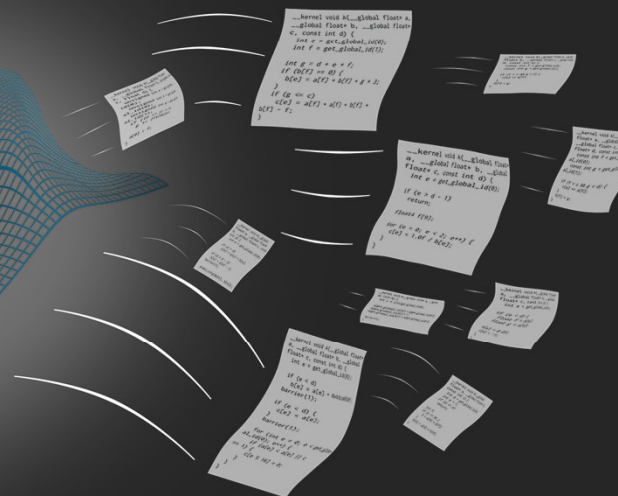
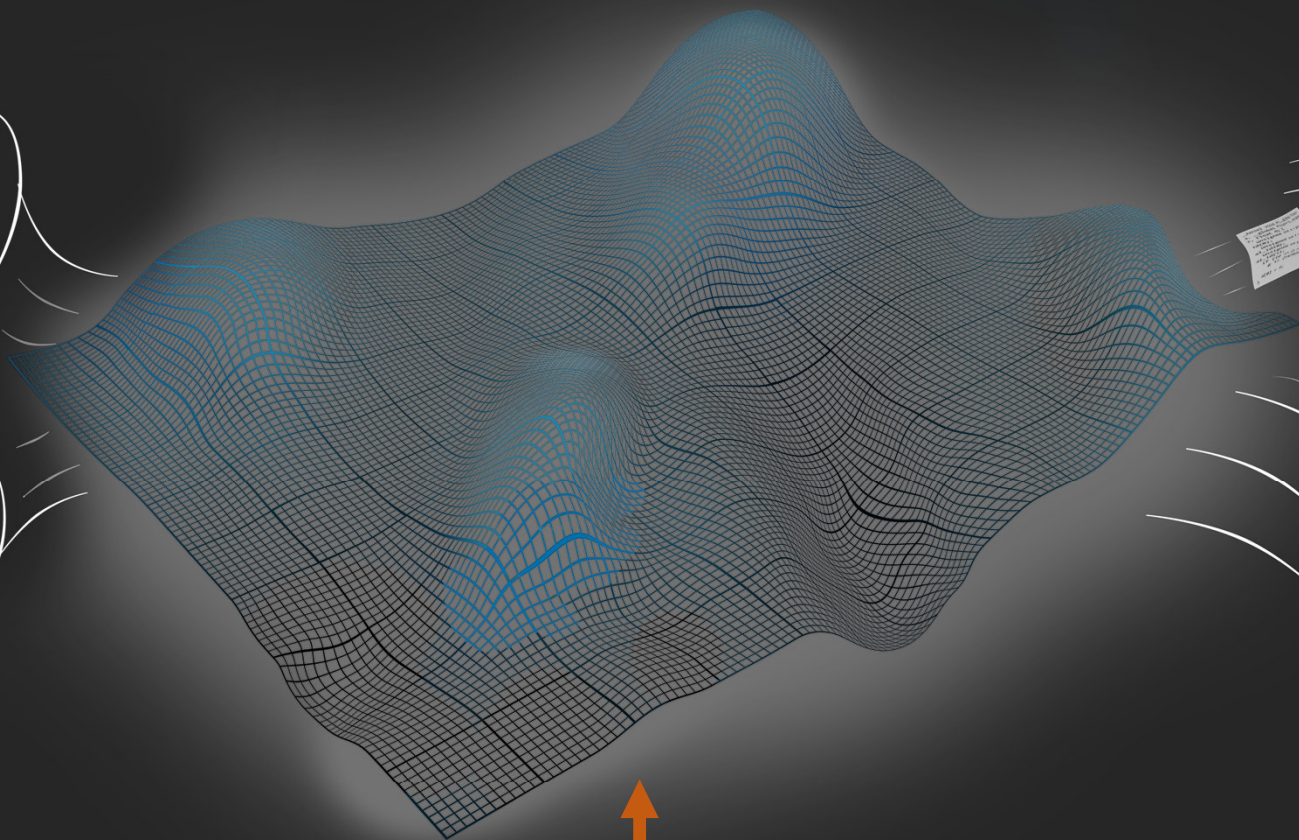
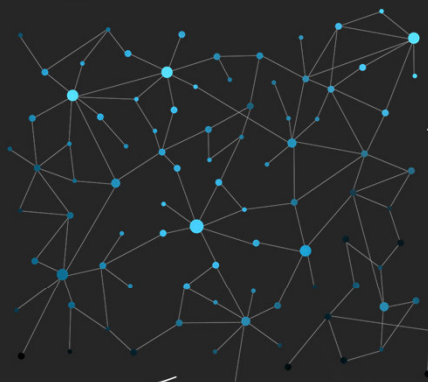
# ***What we need***



**from this**



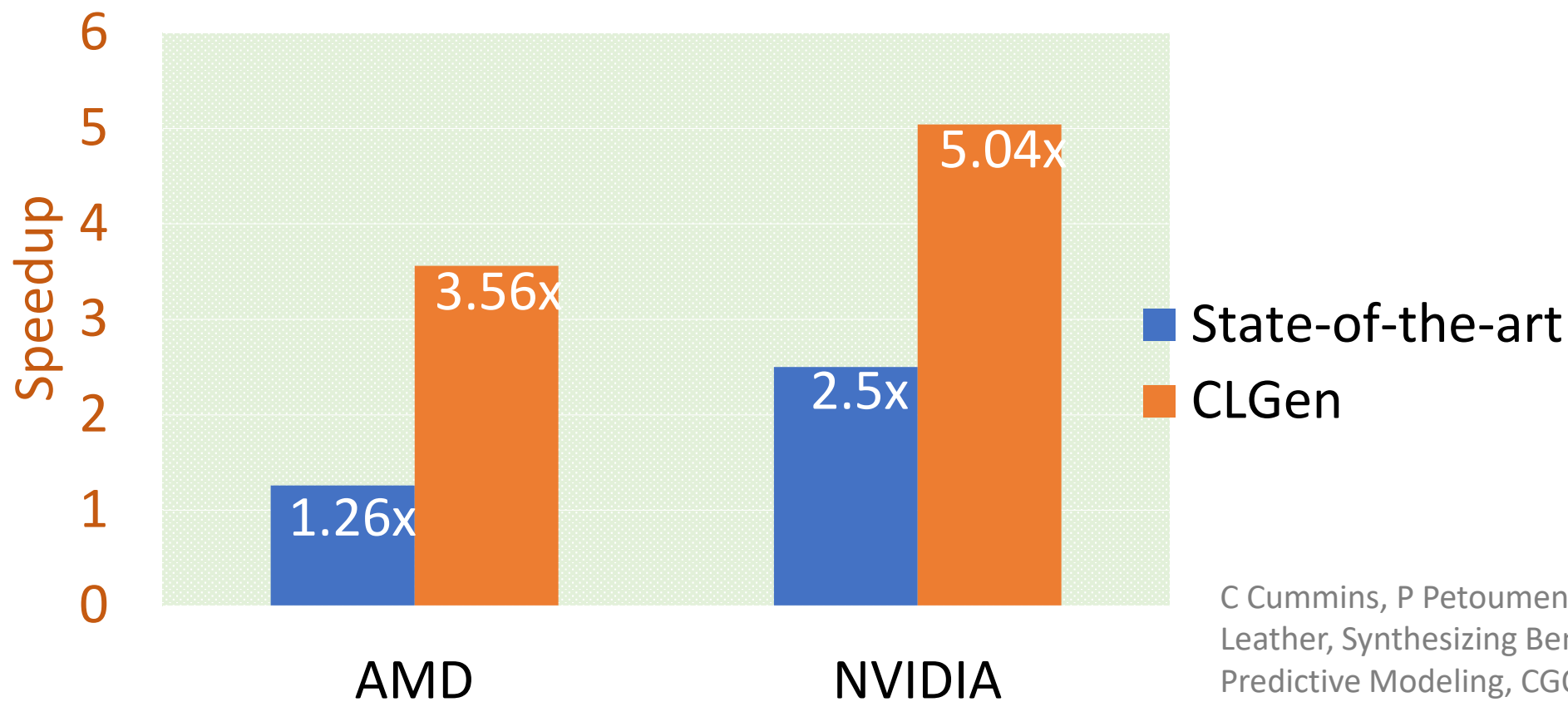
**to this**



model source distr.



# 71 benchmarks, 1,000 synthetic benchmarks



# Language agnostic

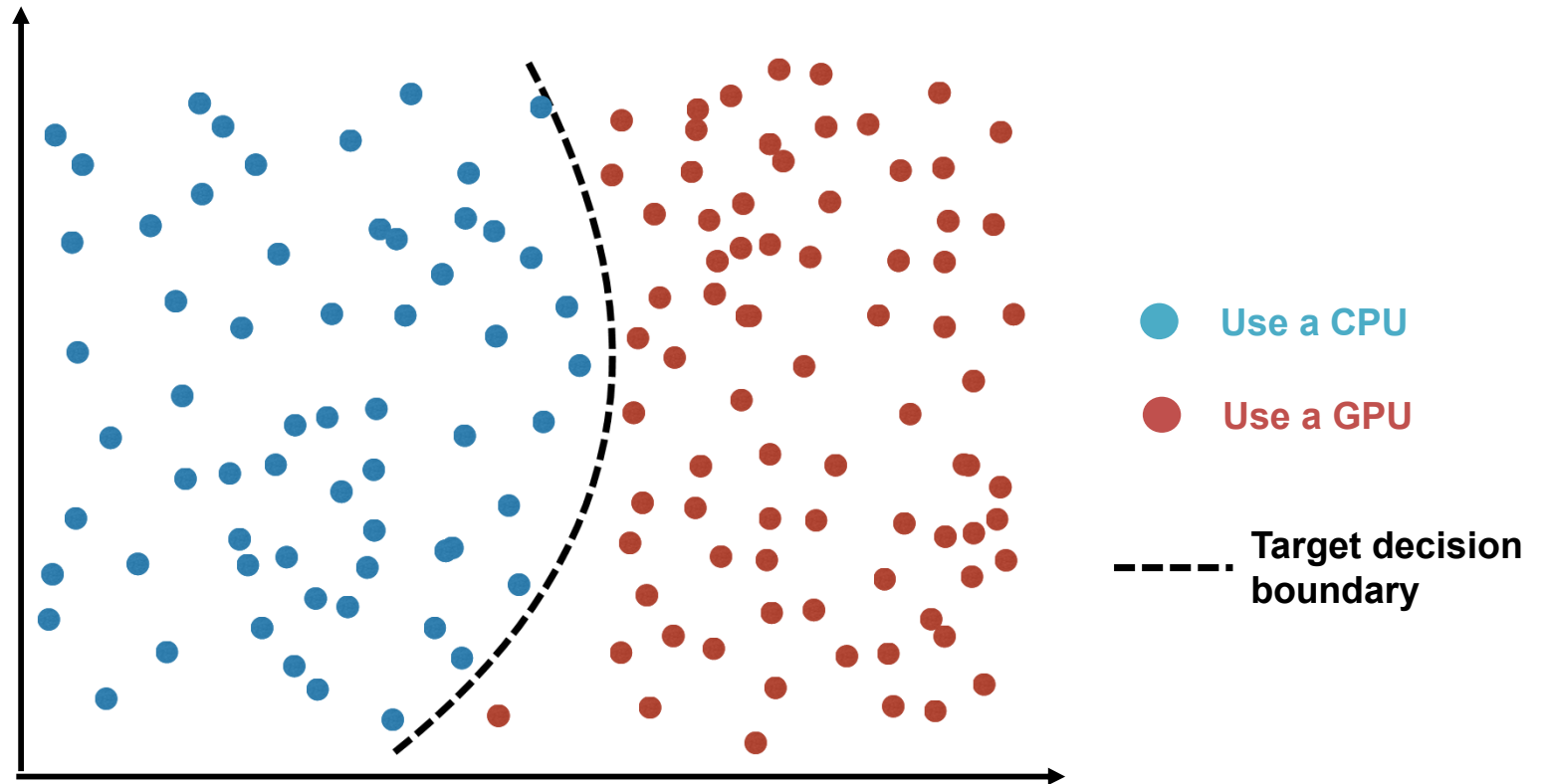
- Generate JavaScript (JS) programs to test JS compilers
- Uncovered **170+ unique** bugs from  
*Chrome V8 , Safari JavaScriptCore, MS edge ChakraCore, Firefox SpiderMonkey, etc.*  
**142 bugs have been verified, 120+ have been fixed.**



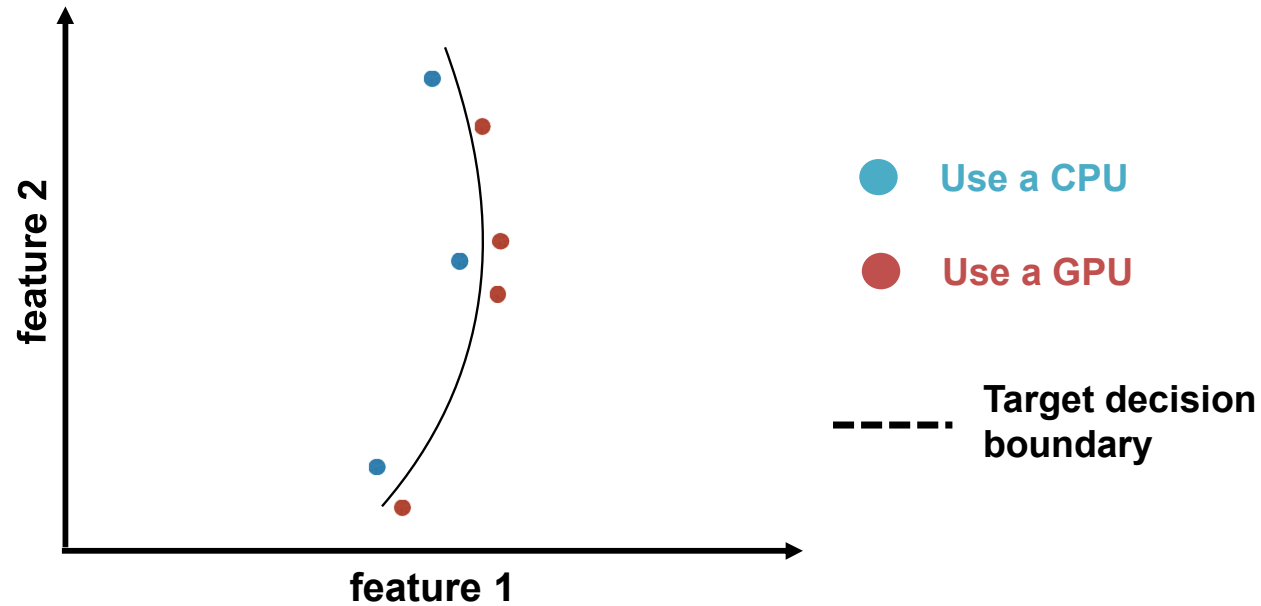
G. Ye, Z. Tang, S. Tan, X. Sun, Z. Wang, *Automated Conformance Testing for JavaScript Engines via Deep Compiler Fuzzing, PLDI 2021*

# ML for predictive modelling– Recap

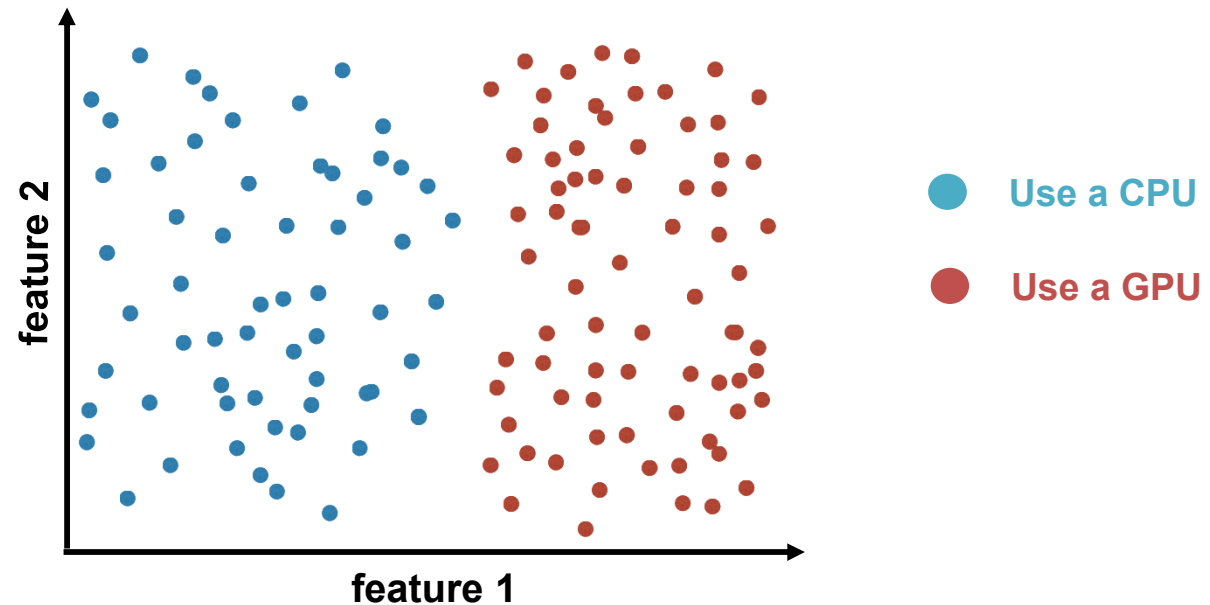
Labelling  
training  
programs can  
be very  
expensive



# Well, we actually only to profile these.



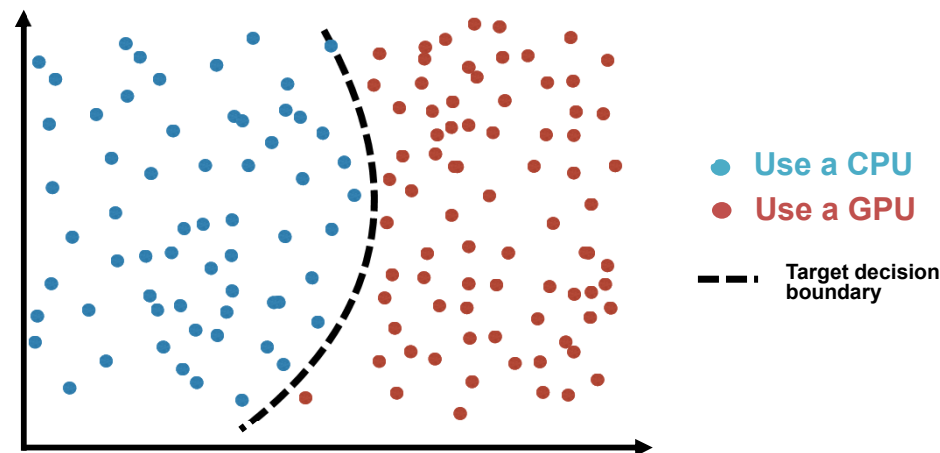
# So these were a complete waste of time!



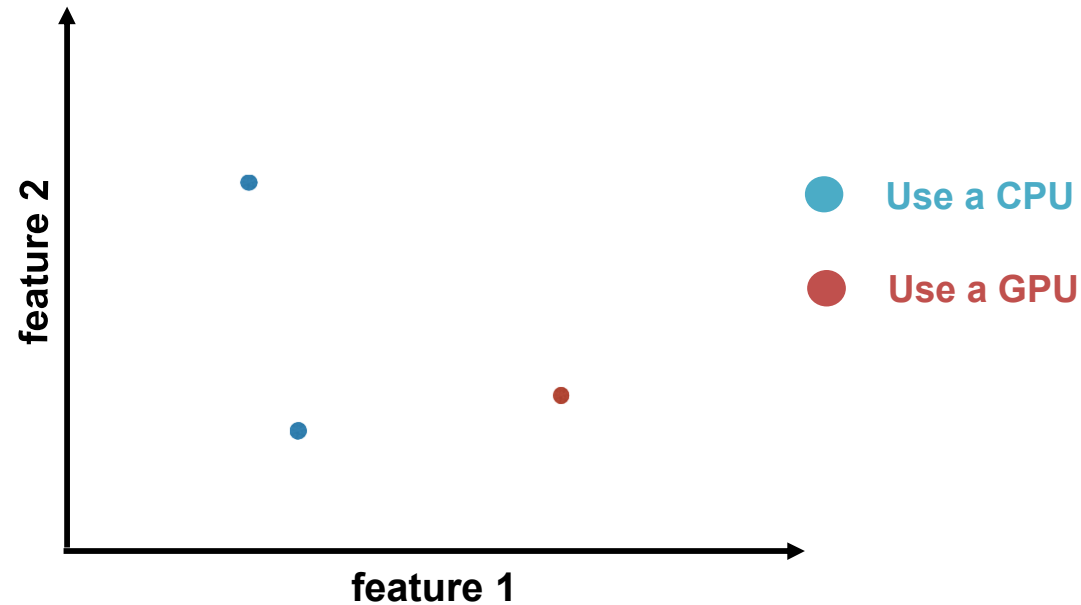
Random profiling inevitably leads to redundancy

# What do we do about it?

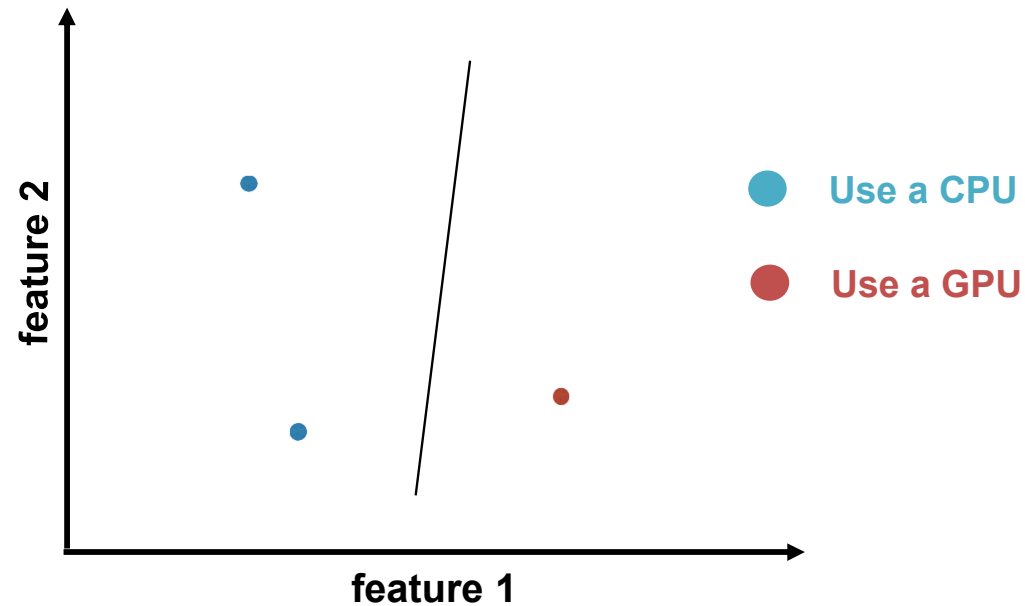
- We cannot know where the informative examples lie
- But, we can let the algorithm make an educated guess – **active learning**



# We start with a few random examples

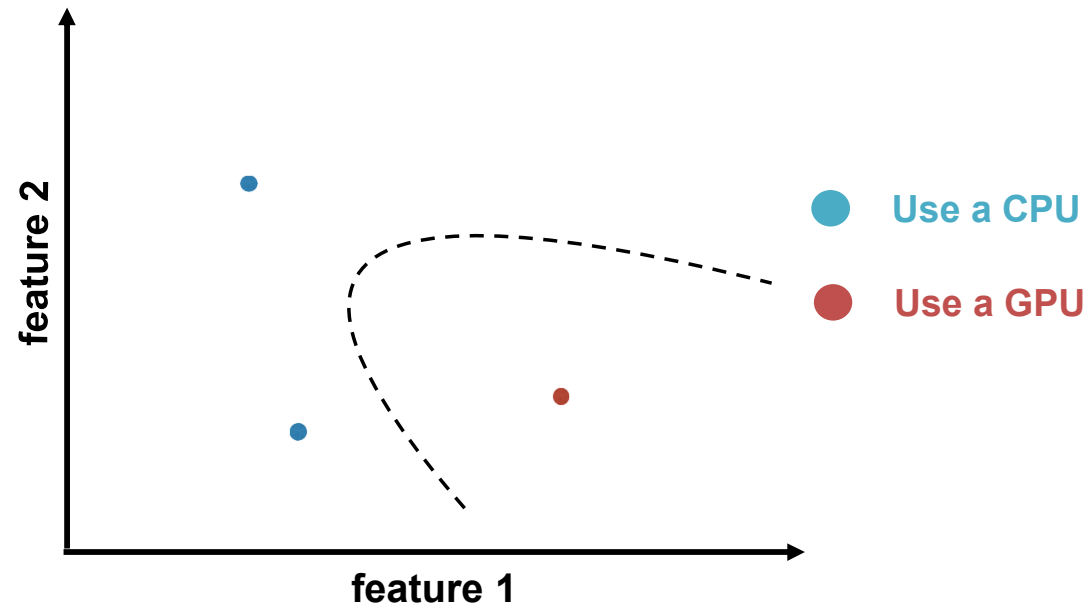


# We form multiple intermediate models

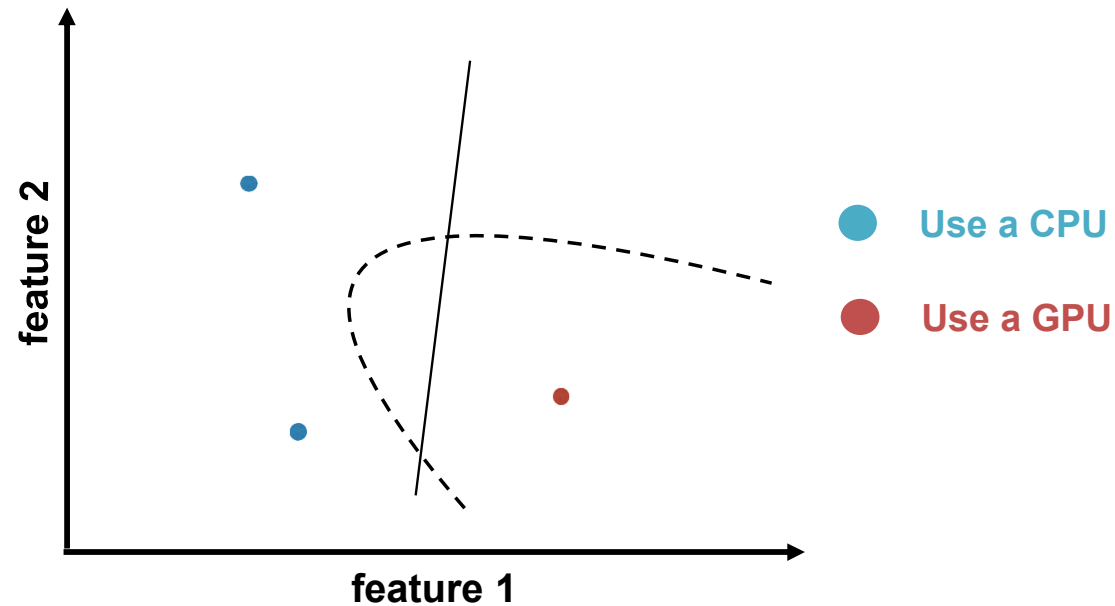




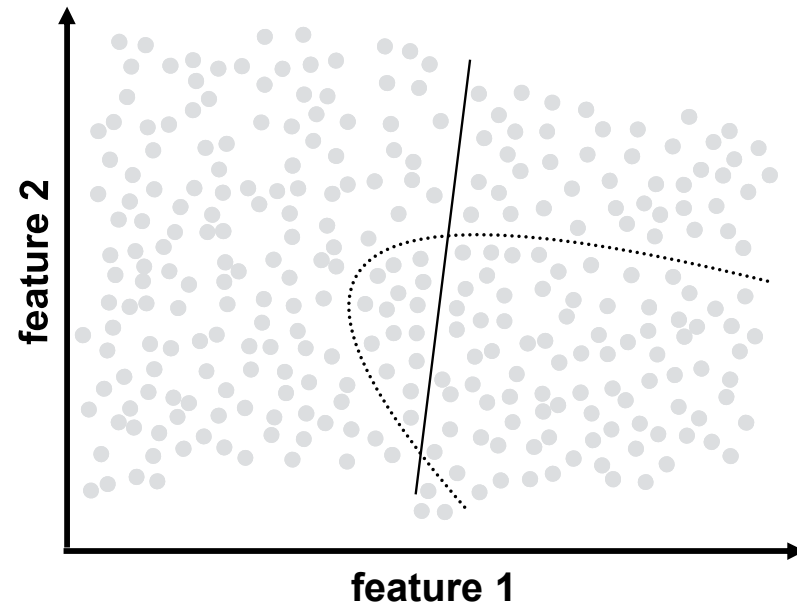
# Each with a distinct algorithm



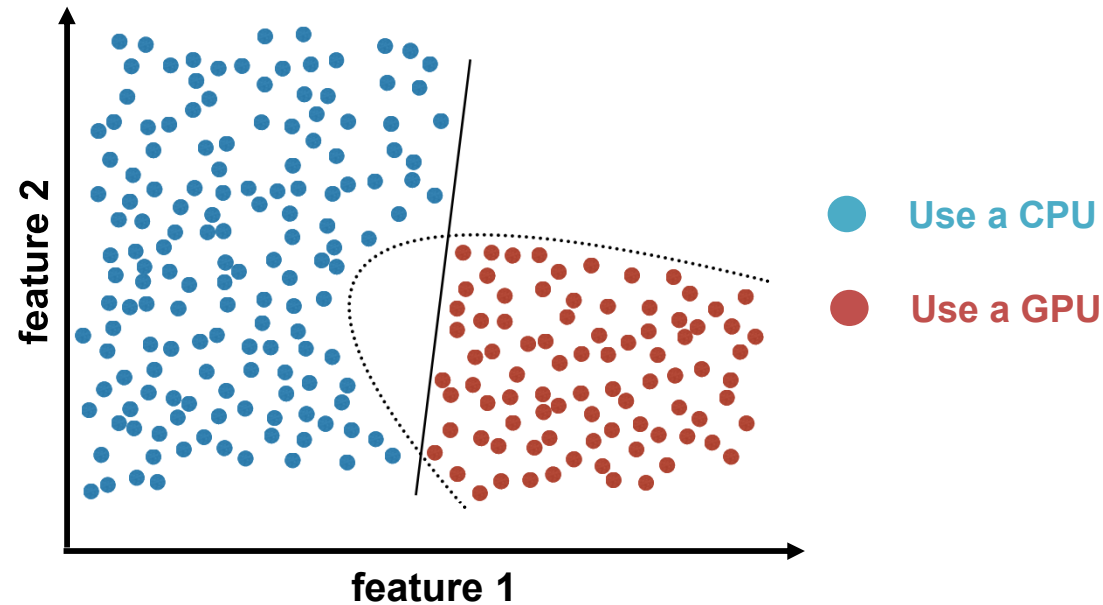
# A “committee” of different models



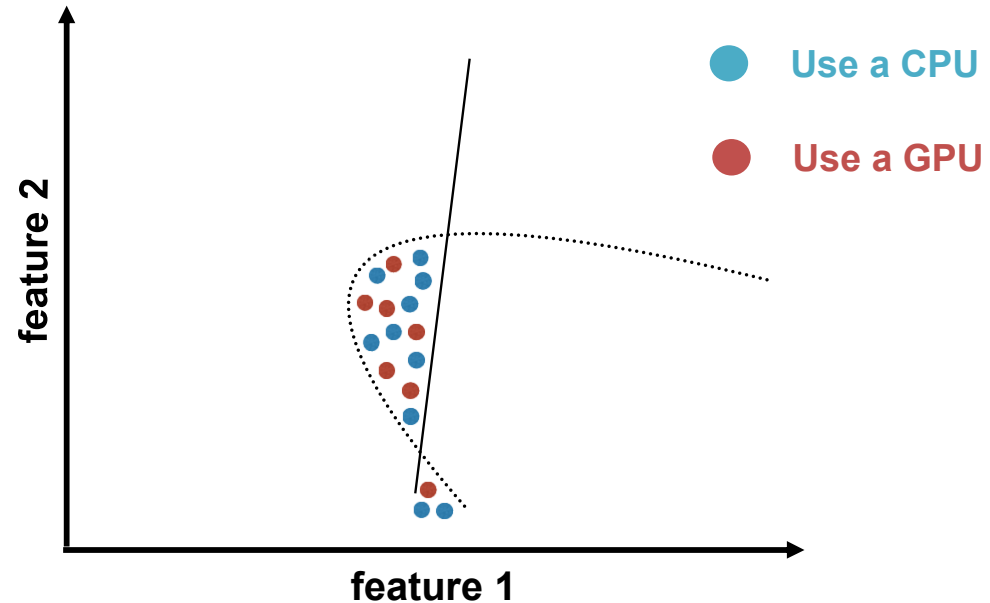
# So what example do we learn from next?



# Broadly the “Committee” will agree...

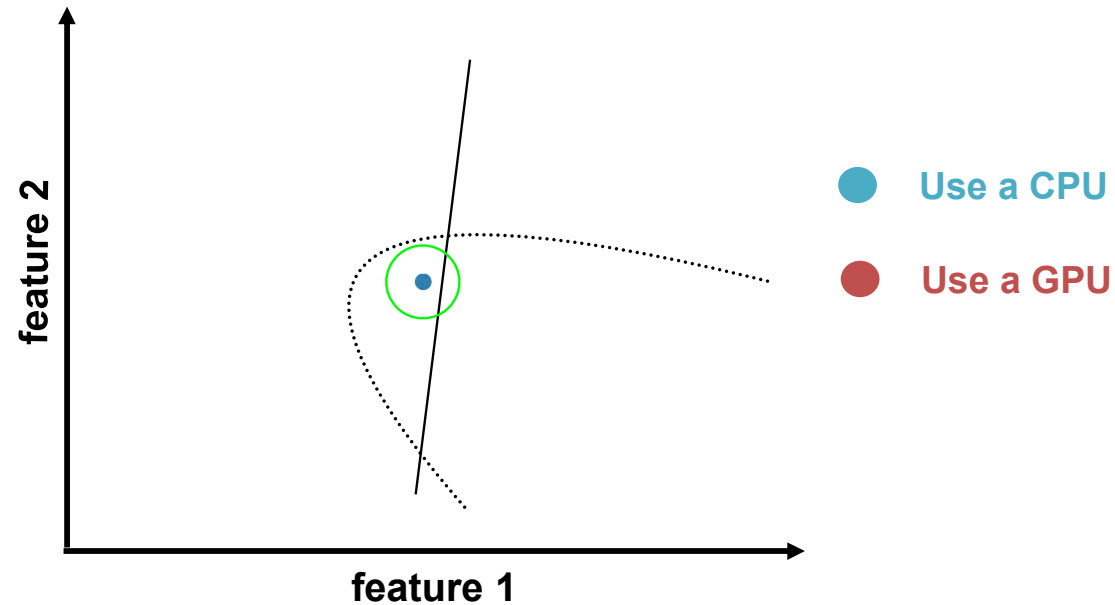


# Sometimes, they disagree...

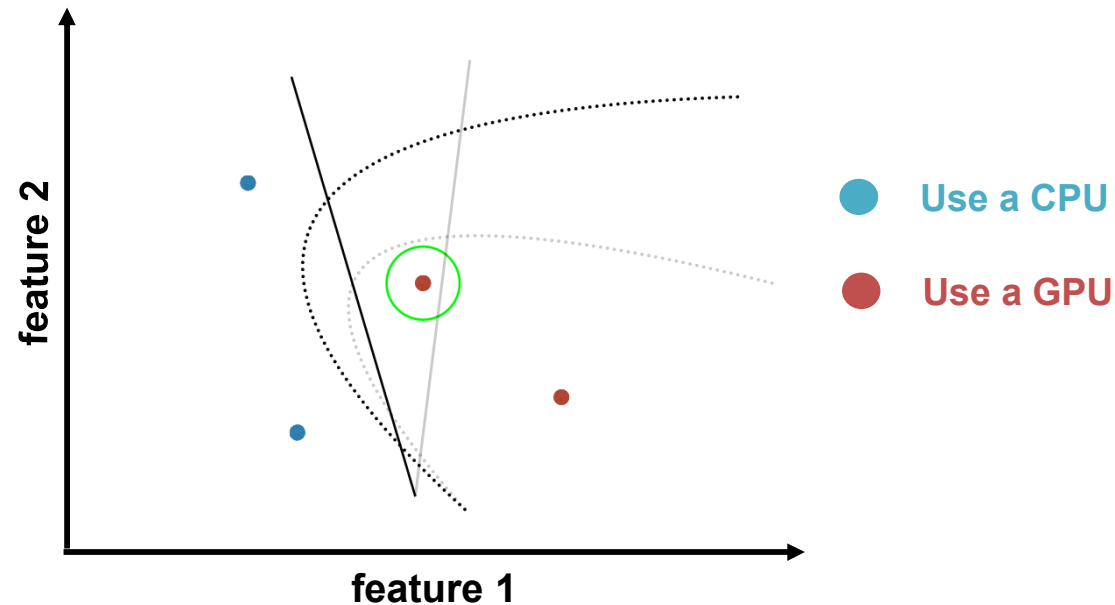


Disagreement regions hold the greatest potential to improve the collective knowledge — learn from there!

# We select one of these examples to label properly



# Then rebuild the intermediate models



Notice the region of disagreement has shrunk  
Eventually the distinct models will converge

# 4x faster on average

## SPAPT Benchmarks

Optimisation space :  $5 \times 10^8$  to  $1 \times 10^{27}$

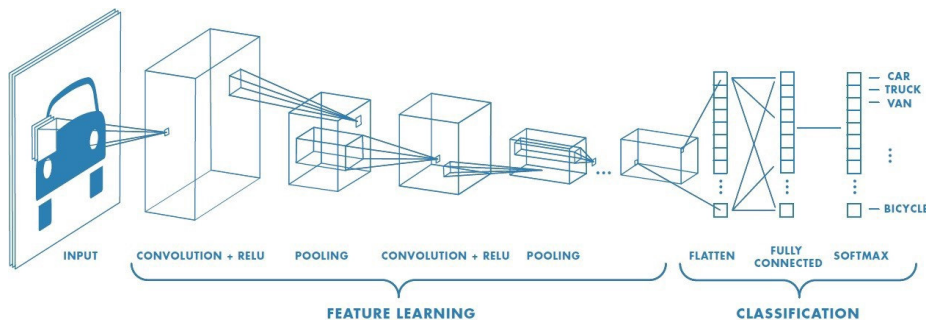
We reduce the profiling time for labelling by 4x.

Minimizing the cost of iterative compilation with active learning, CGO 2017

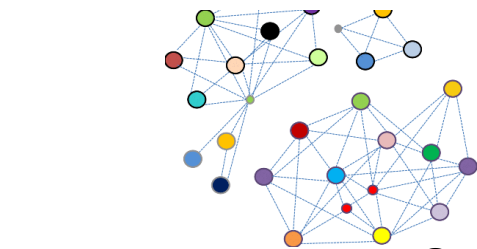


# Challenges ahead

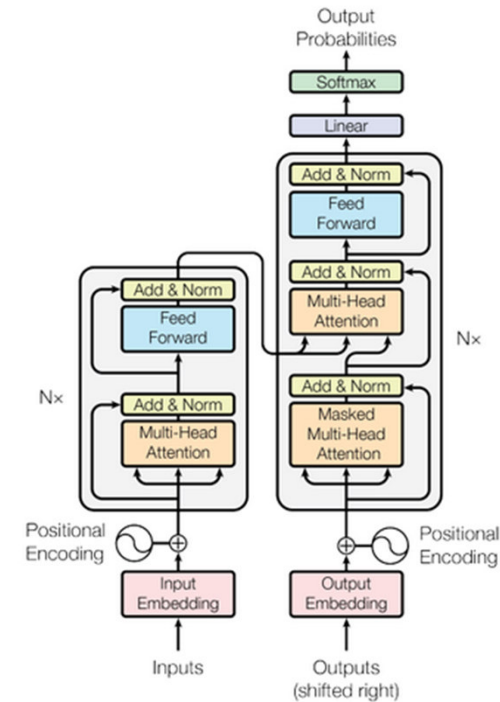
# Representation matters



**Images**  $\Rightarrow$  **Convolutional Neural Networks (CNN)**



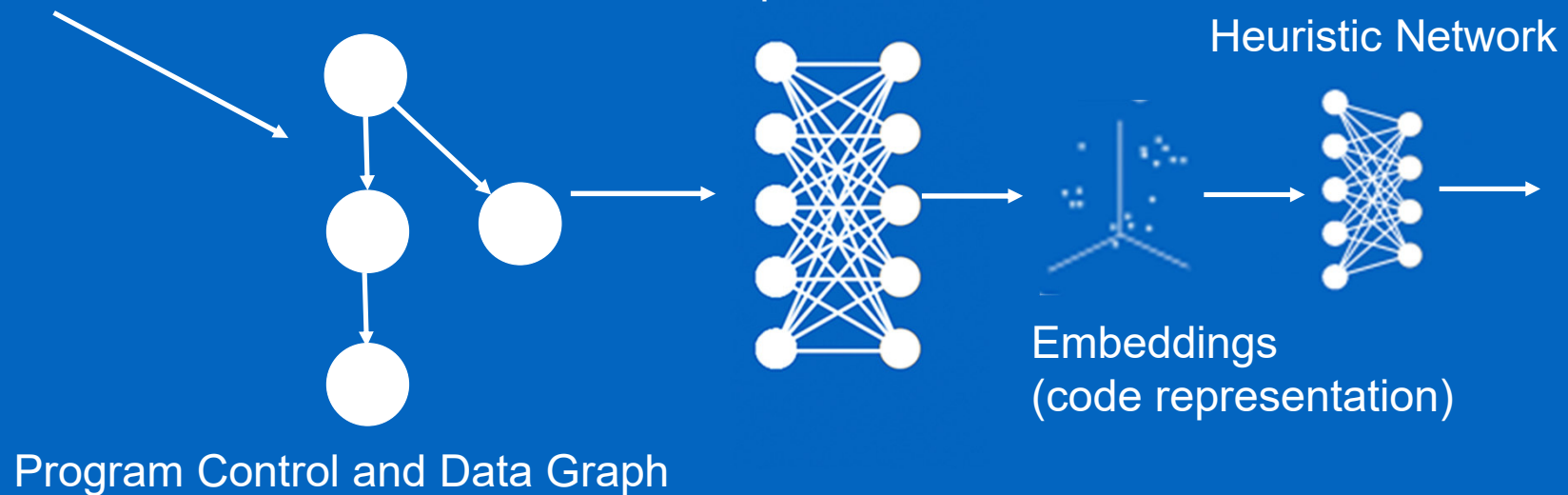
**Social Networks**  $\Rightarrow$  **Graph Convolutional Networks (GCN)**  
**Gated Graph Neural Networks (GGNN)**  
**Graph Transformer Networks (GTN)**



**Text**  $\Rightarrow$  **Transformer based Models**

# How to best represent programs?

```
void memset(void* mem_d, len_t val...)
```



# How to best represent programs?

```
void memset(void* mem_d, len_t val...)
```

Small changes in the program graph (e.g., the loop trip count) can lead to a significant change in the program behaviour.

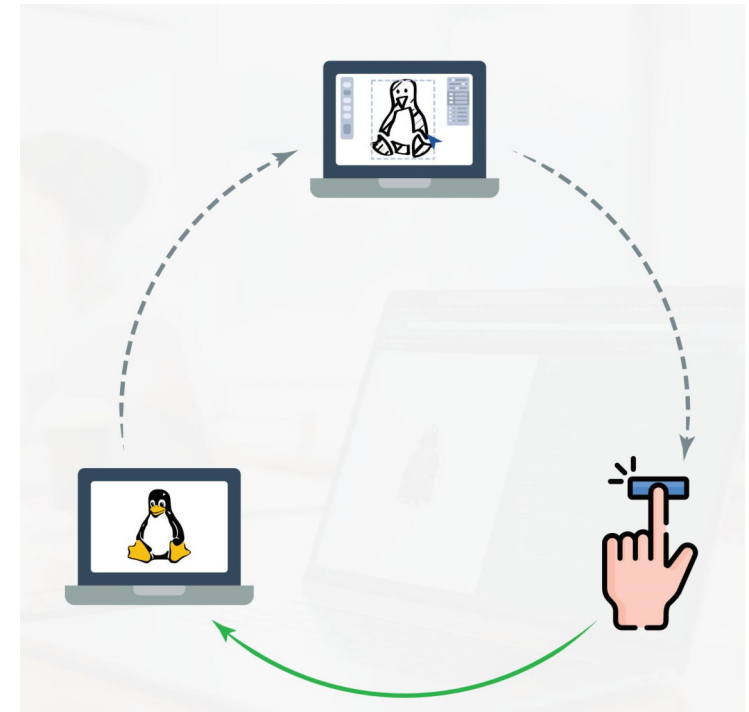
Semantic properties?

# Cloud-based compilation

- Use cloud servers for large-scale compiler analysis and optimisation
  - E.g., pointer alias analysis on large program graphs

Build LLVM in 90 seconds

*From Laptop to Lambda: Outsourcing Everyday Jobs to Thousands of Transient Functional Containers, ATC 2019*



# Conclusions

- **AutoML** to lower entry barriers for ML in compilers
- General-purpose **benchmark synthesis**
- Low-cost profiling with **active learning**
- Many interesting problems ahead

Machine learning in compilers (papers, tools and datasets):  
<https://github.com/zwang4/awesome-machine-learning-in-compilers>

