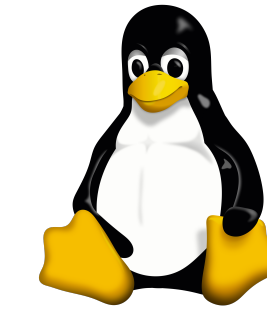


**Building an AI  
that Codes**

<http://chriscummins.cc>

2013



2014



2015



2016



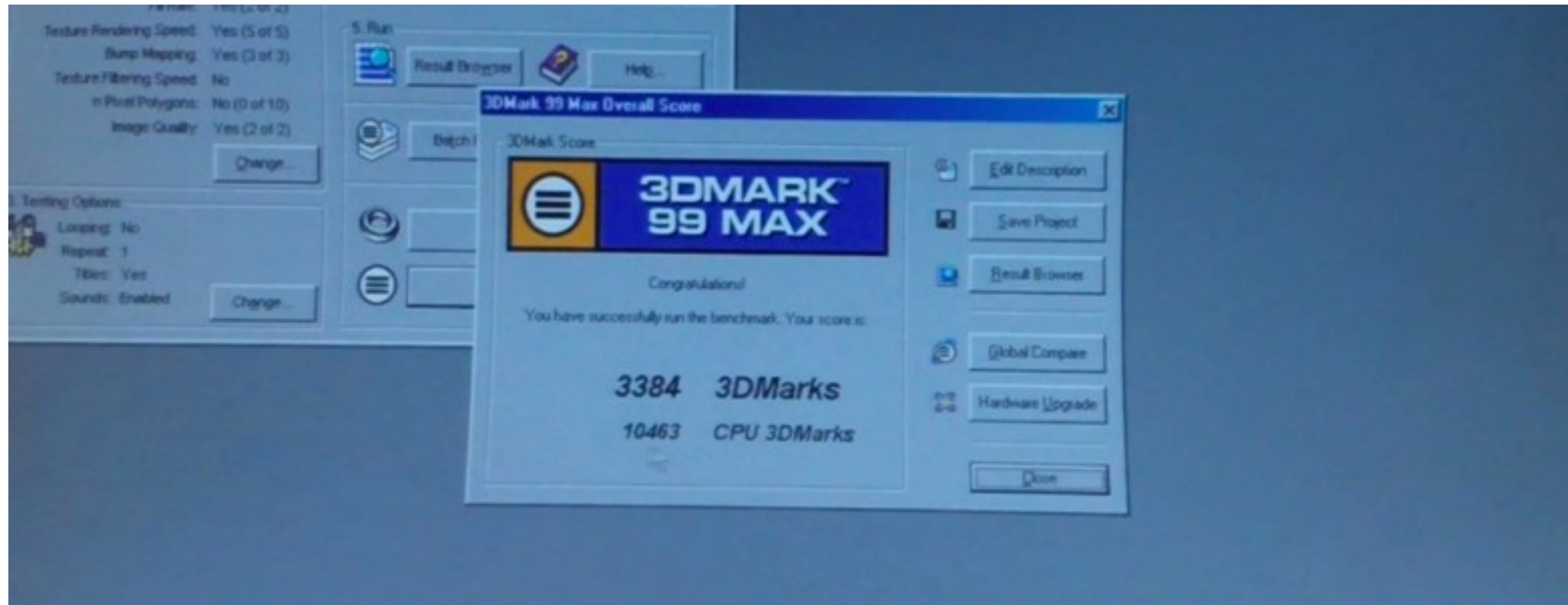








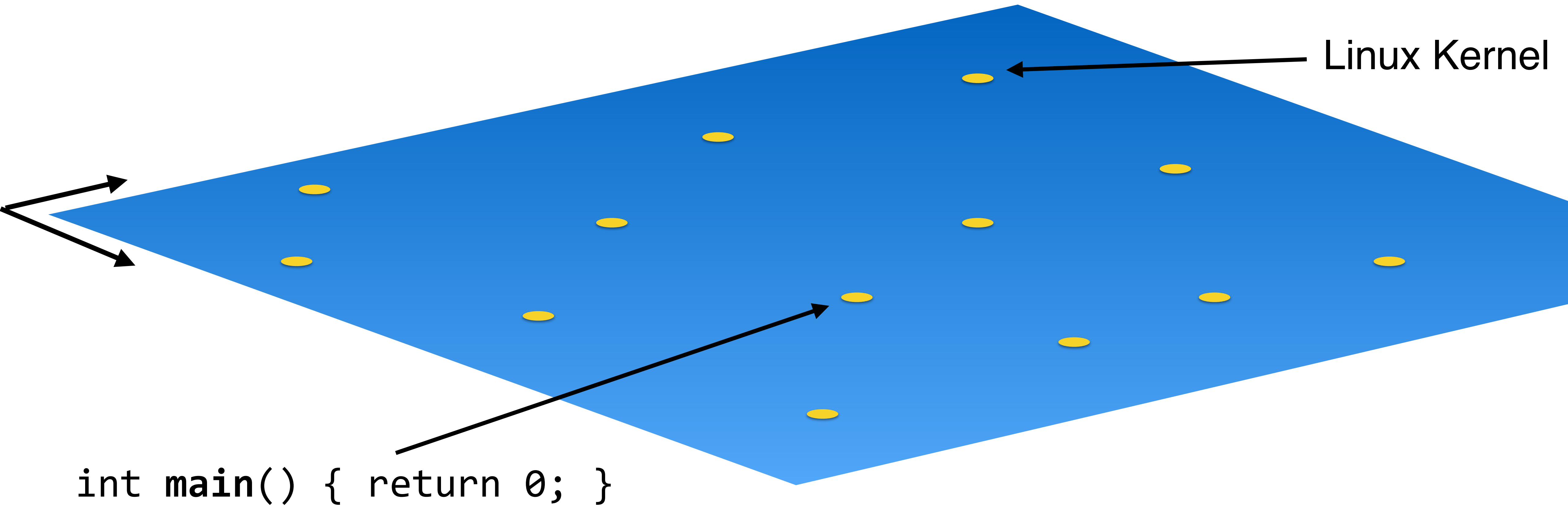
# What makes a good computer?



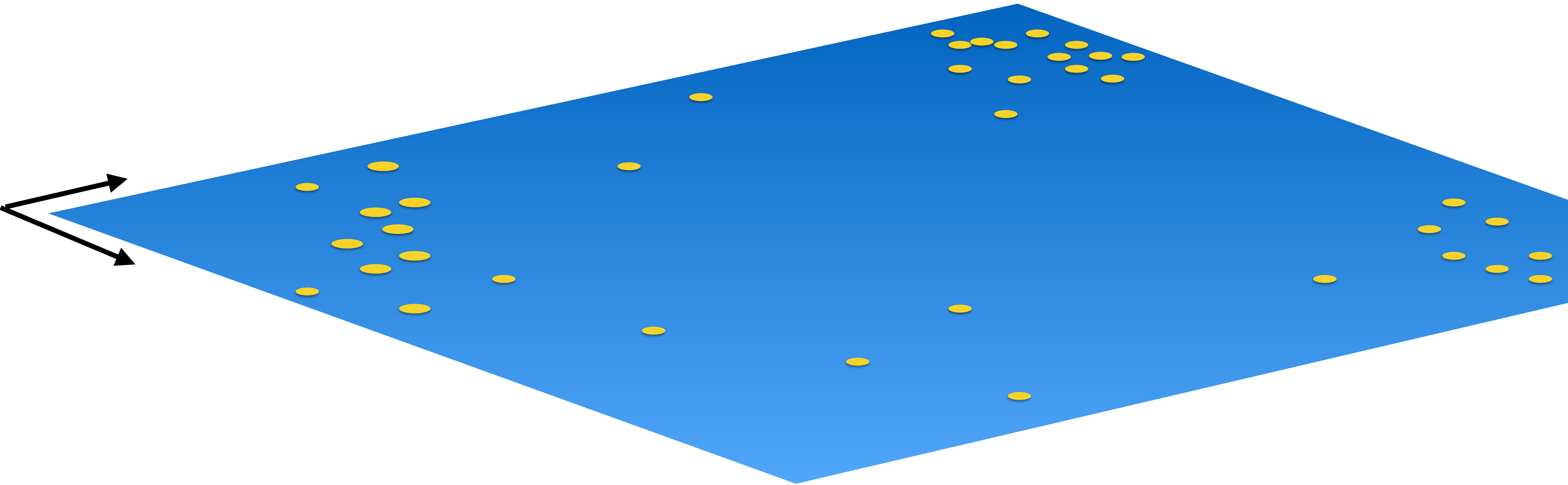
big numbers != smooth games

fast forward ...

Consider the “implementation space”

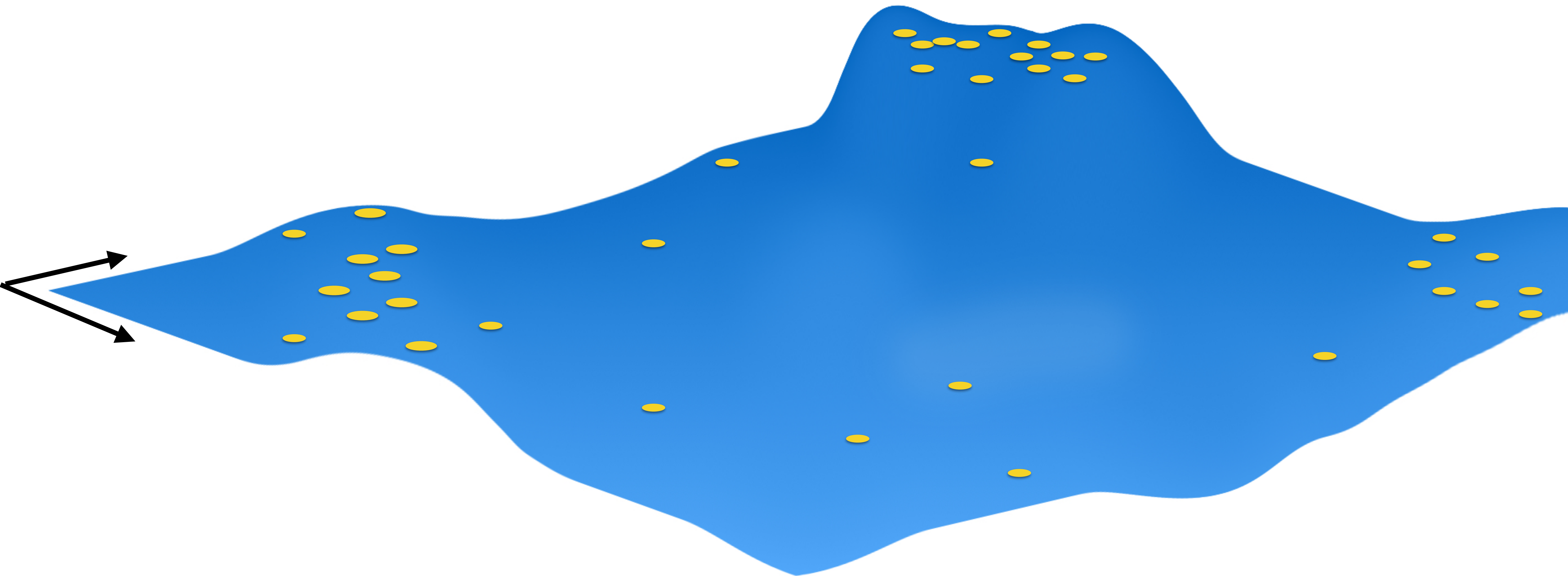


Hypothesis: real source codes form clusters

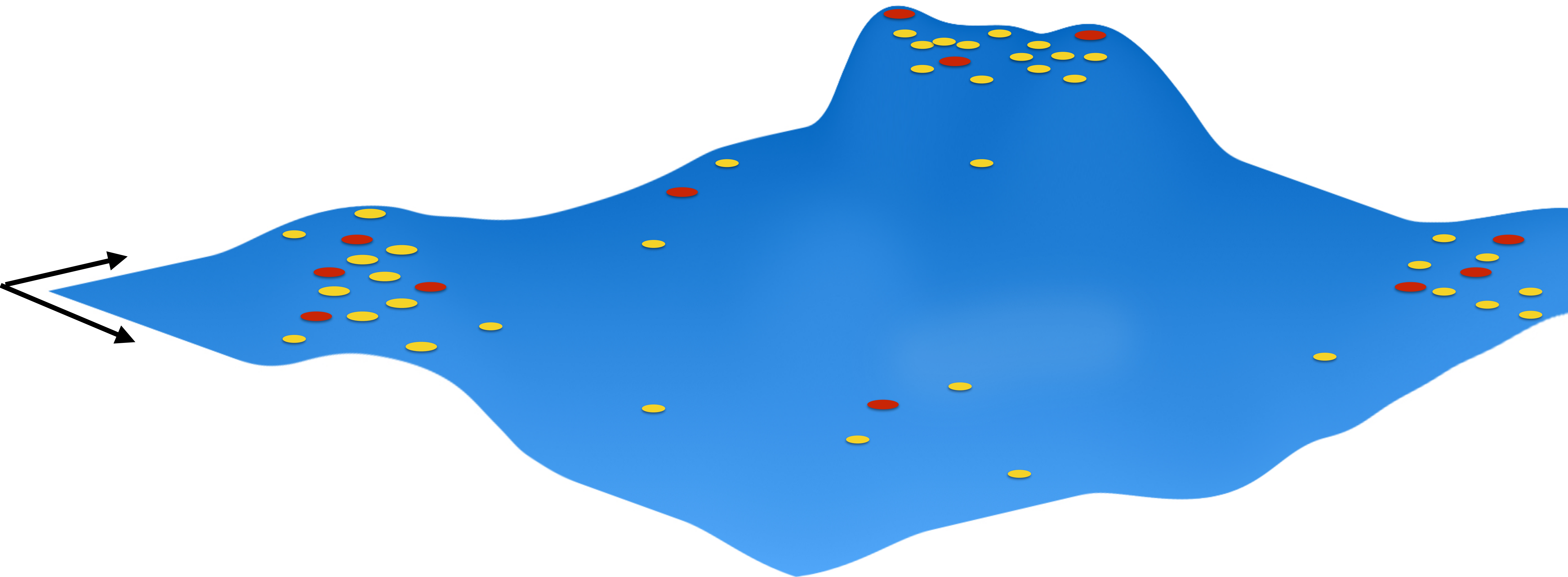




Weight space to match clustering



Sample from weighted space to generate new, representative benchmarks on-demand.



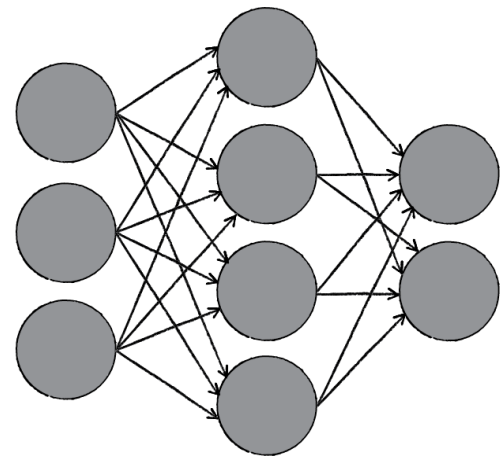
# The idea



Mine programs from



Apply



over implementation space

Generate representative benchmarks on-demand



Human or Robot?

humanorrobot.uk/game/?g=opencl&m=rabt#

Chris

humanorrobot.uk

GamesAbout

Round 1

Player: 1000, Robot: 1000

```
__kernel void A(__global float *a, __global float *b, __global int *c) {
    int d = get_global_id(0);
    float16 e = (float16)db* (*c),db* (*c),
                db* (*c),db* (*c),
                db* (*c),db* (*c),
                db* (*c0),db* (*c1),
                db* (*c2),db* (*c3),
                db* (*c2),db* (*c3),
                db* (*c4),db* (*c5));

    float16 f;
    f = cosh(e);
    a[d * (*c) + 0] = f[0];
    a[d * (*c) + 1] = f[1];
    a[d * (*c) + 2] = f[2];
    a[d * (*c) + 3] = f[3];
}
```

```
__kernel void A(__global float *a, __global float *b, __global float *c) {
    float d;
    float e = 0;
    for (int f = 0; f < 1024; f++) e[e * 16 + g] =
        i[e * (*c) + 0];

    barrier(1);

    for (int g = 1; g < c * (g); ++g) {
        for (int h = 0; h < i - 1; h++) {
            g[h] = f[g] * ((h + 1 & 0xf)) << (1 - i) |
                ((i[h] >> (32 - h))));
        }
        h = (h & 0x80f) + (f - h) >> 1;
        i = f + (e + f) / 2;
        b[f] = i;
    }

    return (f & g);
}
```

This is more human-like

This is more human-like

Score

Rounds Played

alpha

bravo

charlie

delta

© 2016 Chris Cummins.

http://humanorrobot.uk

Curing the Benchmark Deficit: On-Demand Compute Kernel Synthesis using Deep Learning

Chris Cummins

University of Edinburgh

c.cummins@ed.ac.uk

Pavlos Petoumenos

University of Edinburgh

ppetoume@inf.ed.ac.uk

Zheng Wang

Lancaster Univeristy

z.wang@lancaster.ac.uk

Hugh Leather

University of Edinburgh

hleather@inf.ed.ac.uk

Abstract

The quality of performance tuning is bound by the quantity and quality of benchmarks used. Too few benchmarks leads to overfitting; non-representative benchmarks lead to invalid predictions. We present a novel methodology for generating OpenCL compute kernels. Given a corpus of example programs, we apply deep learning across the implementation space, learning a language model from which we obtain new kernels through a process of rejection sampling. We demonstrate our approach for a state-of-the-art machine learning OpenCL autotuner. With the addition of synthesised compute kernels, we improve the accuracy of machine learning predictions from XXX% to XXX%, demonstrating up to  $X\times$  speedup over the hand-selected benchmarks.

**Keywords** Synthetic program generation, OpenCL, Deep Learning, GPUs

1. Introduction

Benchmarking parallel applications is hard. State of the practice is lacking [1]. OpenCL benchmark suites: Rodinia [2], Parboil [3], Polybench [4], SHOC [5], AMD SDK<sup>1</sup> and NVIDIA SDK<sup>2</sup>. TODO: Tease the small number of benchmarks used. Benchmarking OpenCL [6]. Are benchmarks suites representative? Exploring the full performance spectrum [7]. Characterising workloads of Rodinia and Parsec [8]. In previous works we used stochastic template substitution to generate stencil benchmarks for autotuning [9, 10]. This template based approach is not general

2. Motivation

We surveyed the benchmarking methodologies of GPU research papers from top tier conferences between 2013–2016: CGO, HiPC, PACT, and PPOPP. By aggregating the sources and quantities of benchmark kernels from 27 papers, we find that 79% of benchmark kernels come from four benchmark suites<sup>3</sup>. Figure 1 and Figure 2.

3. Generating Compute Kernels

Overview of methodology. Figure 3.

purpose, requires laborious human effort, and does not guarantee to generate representative benchmarks.

*RQ1: Can the quality of machine learning predictions be improved with the addition of representative benchmarks?*

For this, we need a methodology for generating novel source codes. This leads to the subsequent research question:

*RQ2: Given a program checker and a corpus of example programs, can a language model learn to generate new programs?*

<sup>1</sup><http://developer.amd.com/>

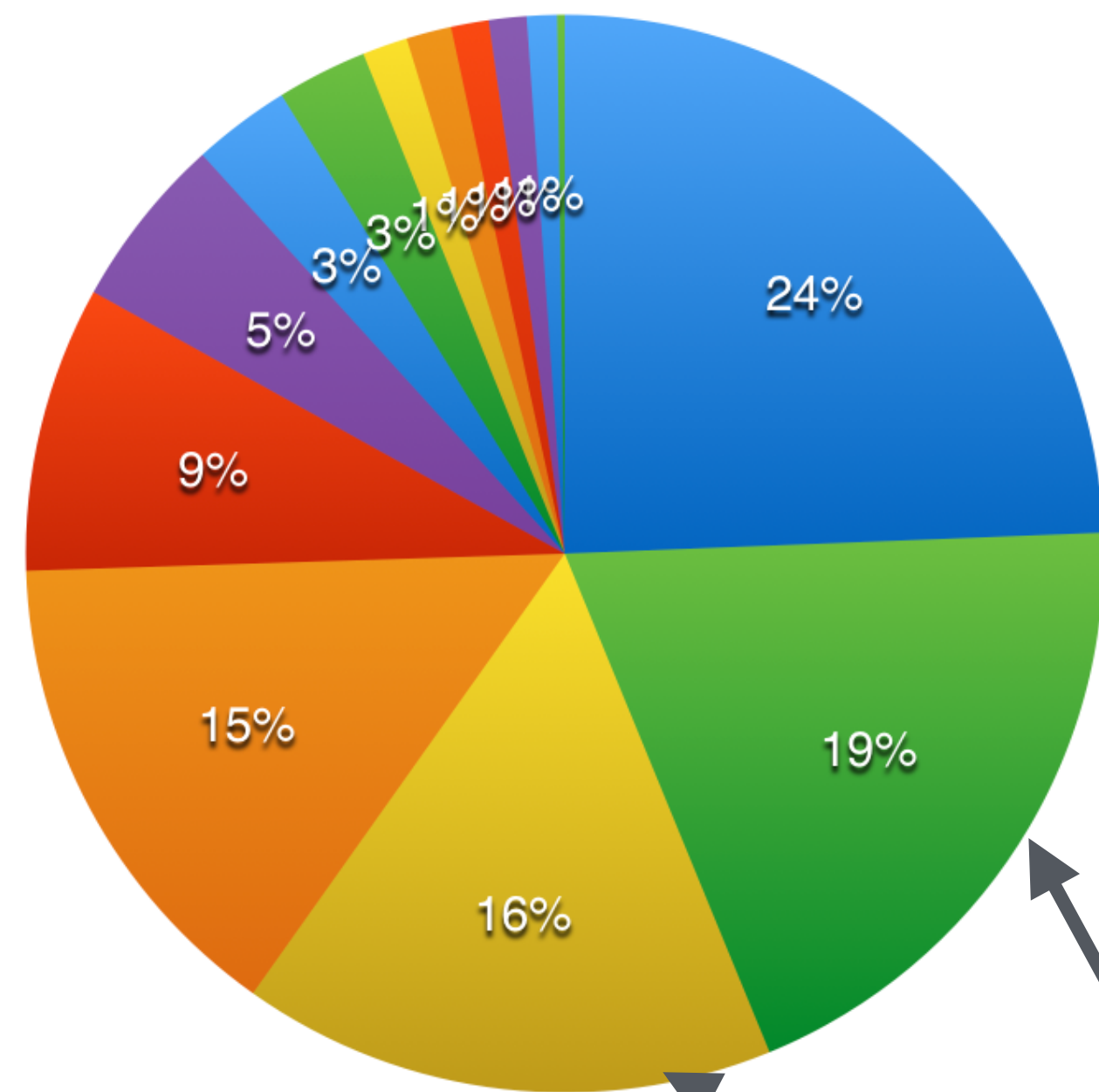
<sup>2</sup><https://developer.nvidia.com/opencl/>

<sup>3</sup>Raw data available at: <http://bit.ly/foo>

1

2016/7/18

2013-2016 state of practise  
27 top-tier GPU papers



79% of results from **4**  
benchmark suites

(SDK sample codes)