

Approximating Computer Architecture

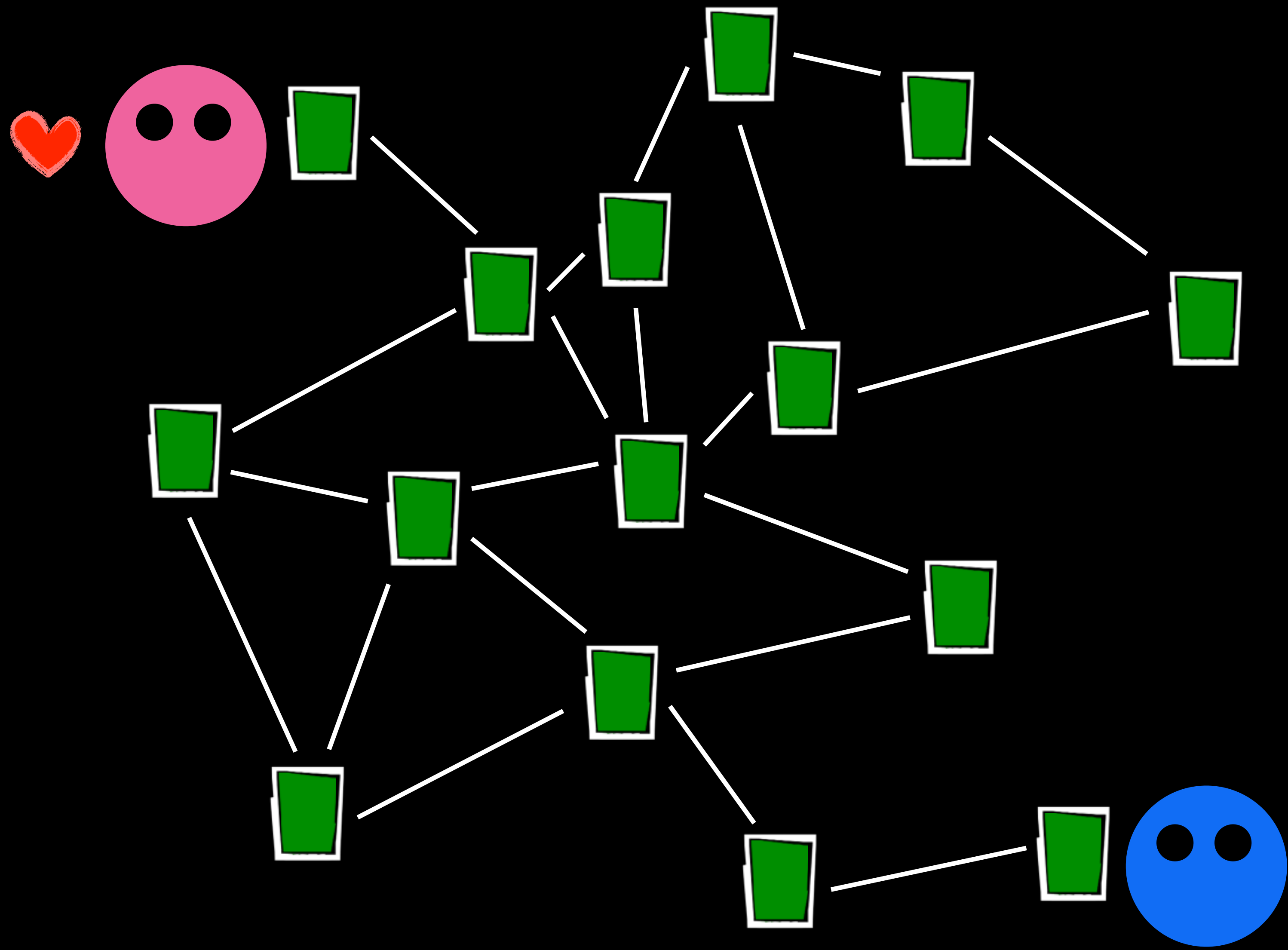
Adrian Sampson Cornell

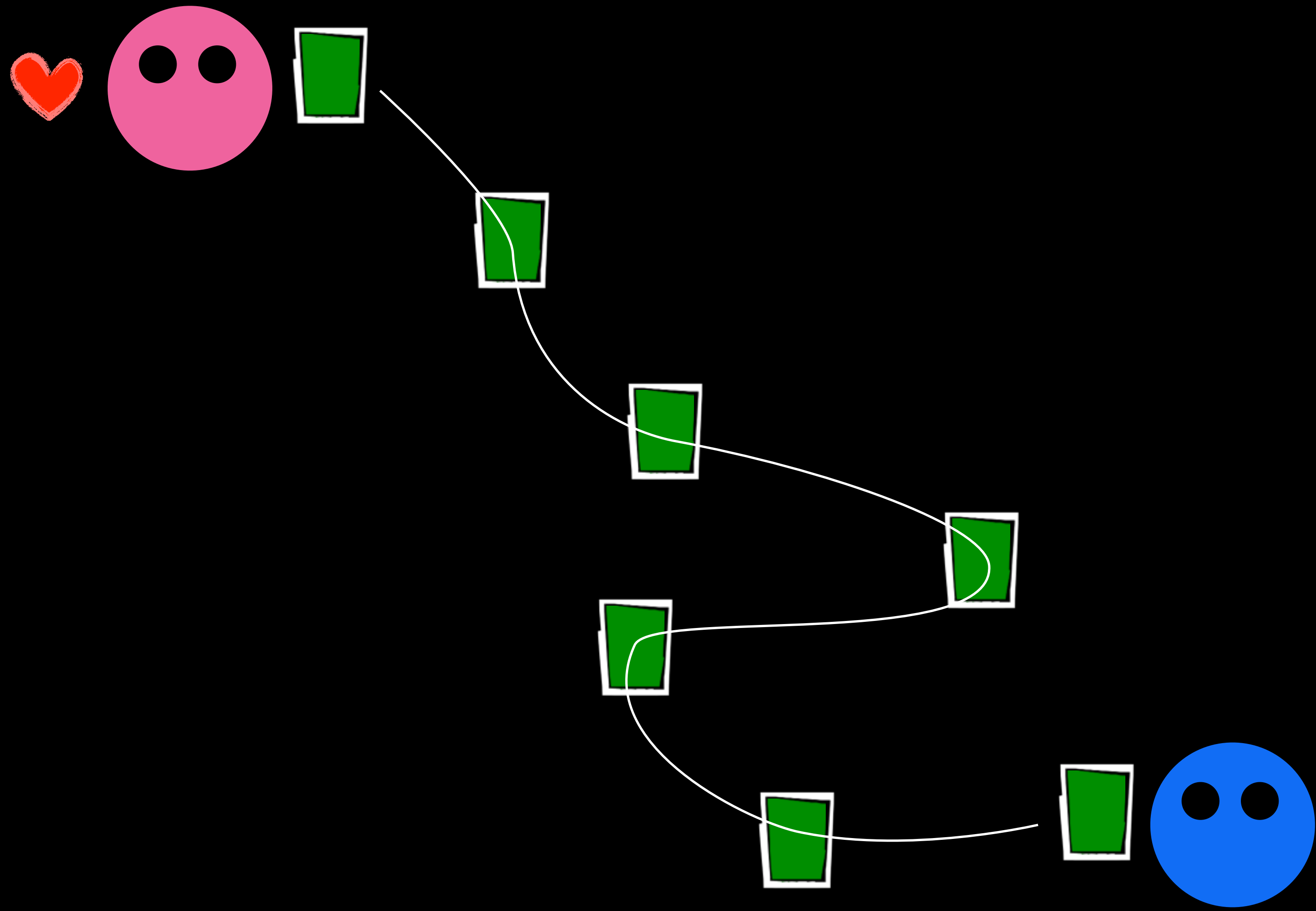
SOME DUBIOUS ADVICE!

1. NEVER LISTEN TO ADVICE
2. ALWAYS LISTEN TO ADVICE
3. IT'S NOT NAIVETÉ; IT'S BEGINNER'S MIND
4. THE BEST COMPUTER ARCHITECTURE RESEARCH ALWAYS "CHEATS"
5. DON'T DO FRAMEWORKS
6. START A BLOG
7. TALK TO YOUR FRIENDS

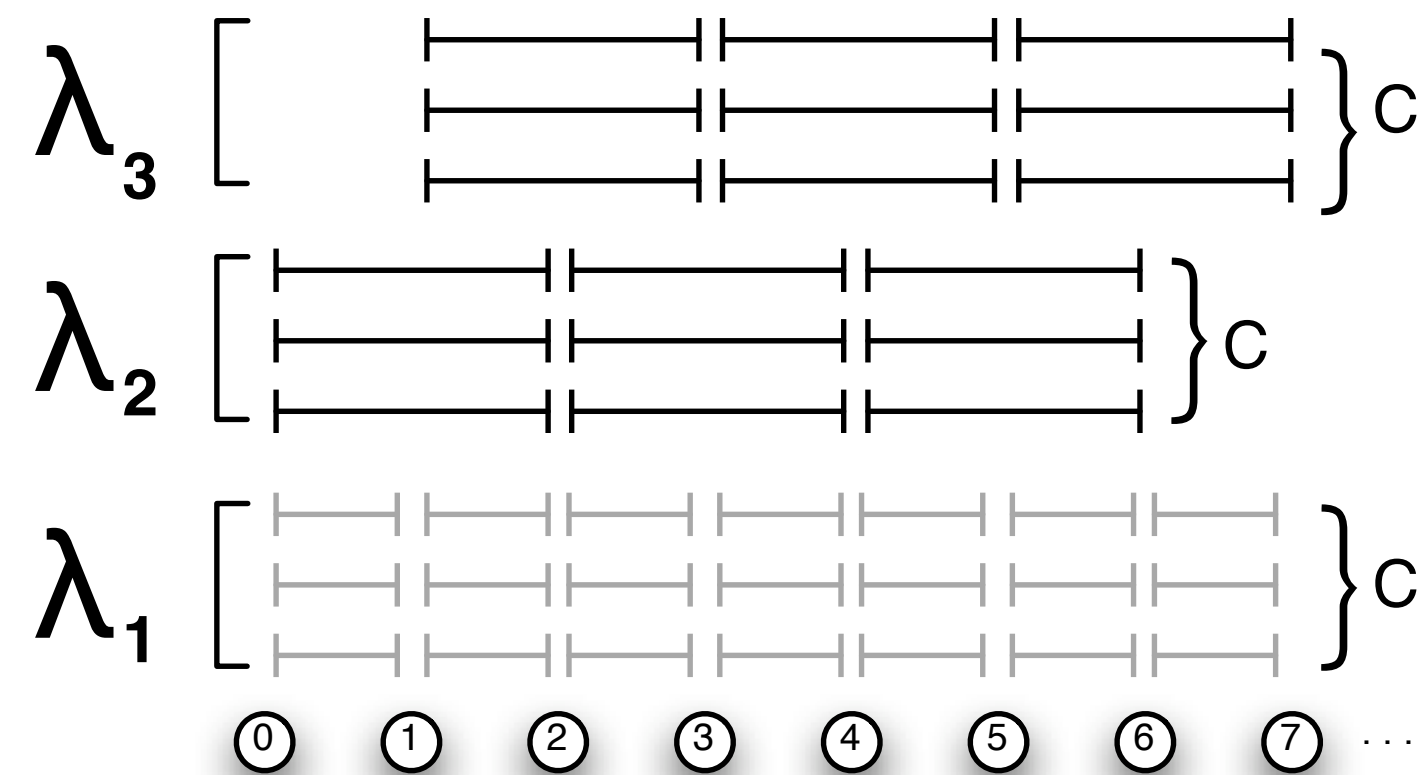


HARVEY
MUDD
COLLEGE
(YESTERDAY!)





virtual topology



traffic grooming

Dynamic Grooming Algorithm (DGA)

Input: A connection request \mathbf{c} and a network state function U .

Output: Returns *true* if \mathbf{c} was satisfied and *false* otherwise. If it was satisfied, U and $S(\mathbf{c})$ reflect the new network state.

```

1  $S(\mathbf{c}) \leftarrow \emptyset$ 
2  $s \leftarrow (L(\mathbf{c}), \min\{R(\mathbf{c}), L(\mathbf{c}) + r\})$ 
3 while  $L(s) \neq R(s)$  do
4   if  $U(s) < C$  then
5      $U(s)++$ 
6     let  $S(\mathbf{c})$  contain  $s$ 
7      $s \leftarrow (R(s), \min\{R(\mathbf{c}), R(s) + r\})$ 
8   else
9     // a conflict occurs with  $s$ 
10     $s \leftarrow (L(s), R(s) - 1)$ 
11  end
12 return  $R(s) == R(\mathbf{c})$  // true iff connection
    is successfully satisfied

```



← RAN
LIBESKIND -
HADAS,
HARVEY MUDD

1. NEVER LISTEN TO ADVICE
2. ALWAYS LISTEN TO ADVICE
3. IT'S NOT NAÏVETÉ; IT'S BEGINNER'S MIND
4. THE BEST COMPUTER ARCHITECTURE RESEARCH ALWAYS "CHEATS"
5. DON'T DO FRAMEWORKS
6. START A BLOG
7. TALK TO YOUR FRIENDS

UNIVERSITY OF
WASHINGTON

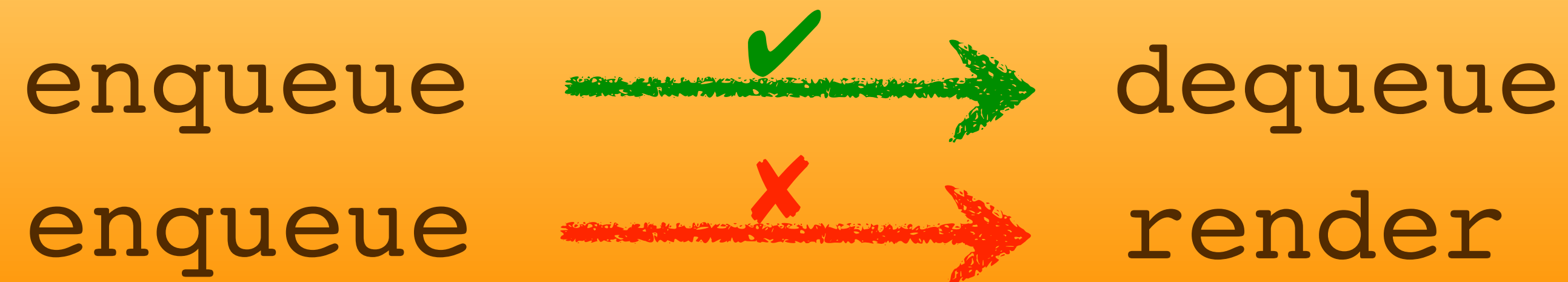


Code-Communication Specifications

Writer Thread

Reader Thread

What **code** may communicate across threads?



May writes in enqueue be read by other threads in dequeue?

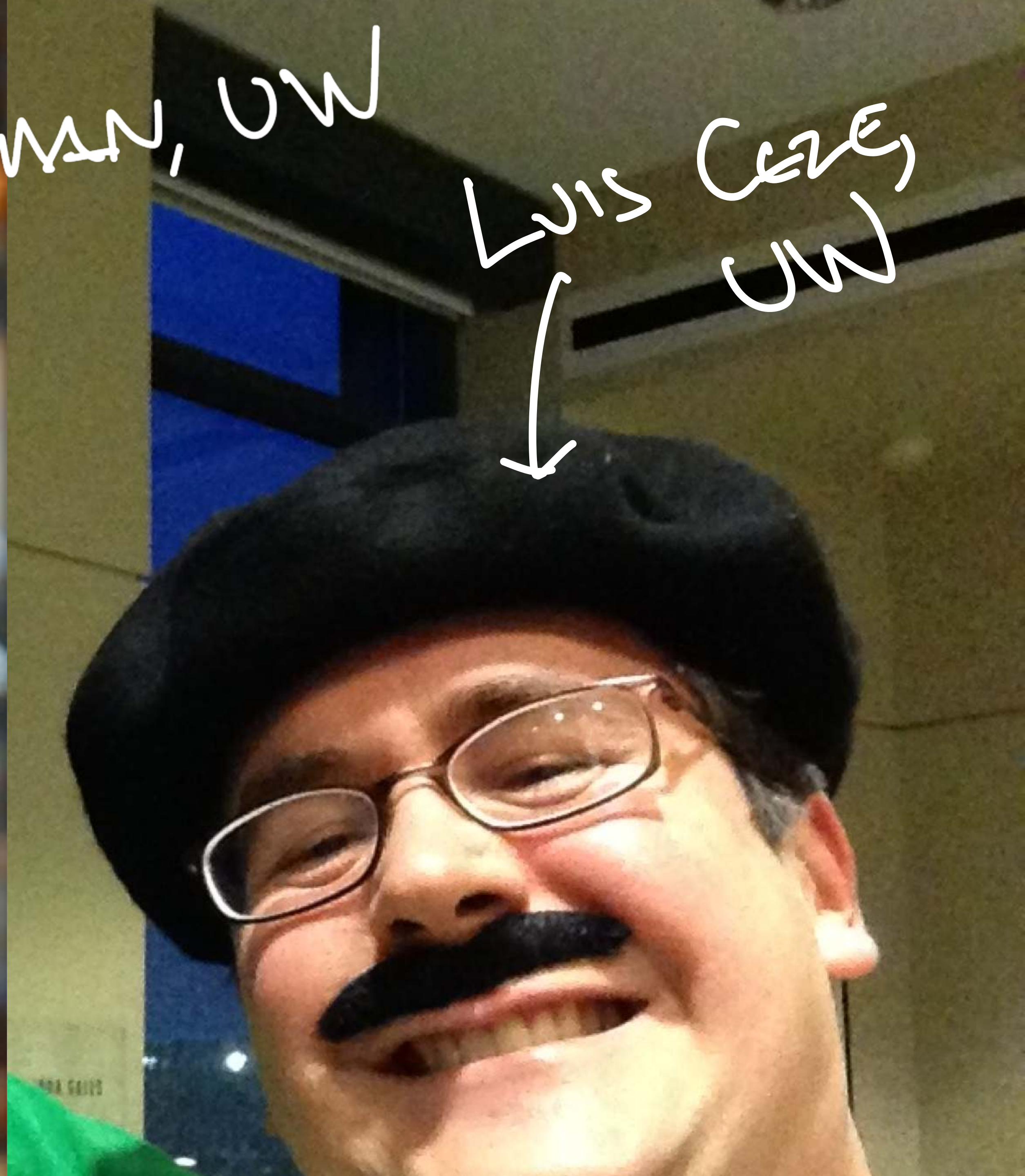


Is this ~~unim~~?

COMPUTER ARCHITECTURE?



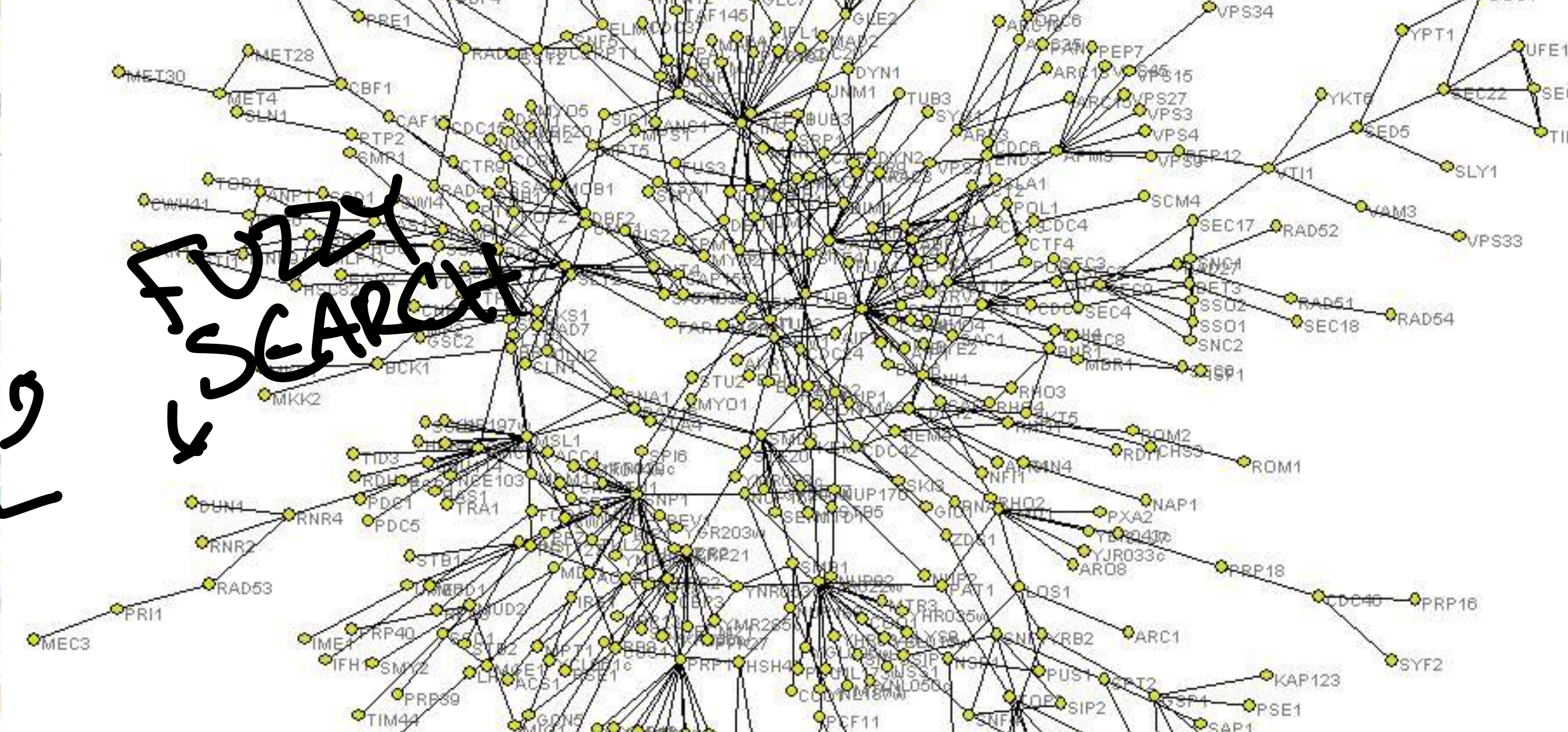
DAN
CROSSMAN, UW



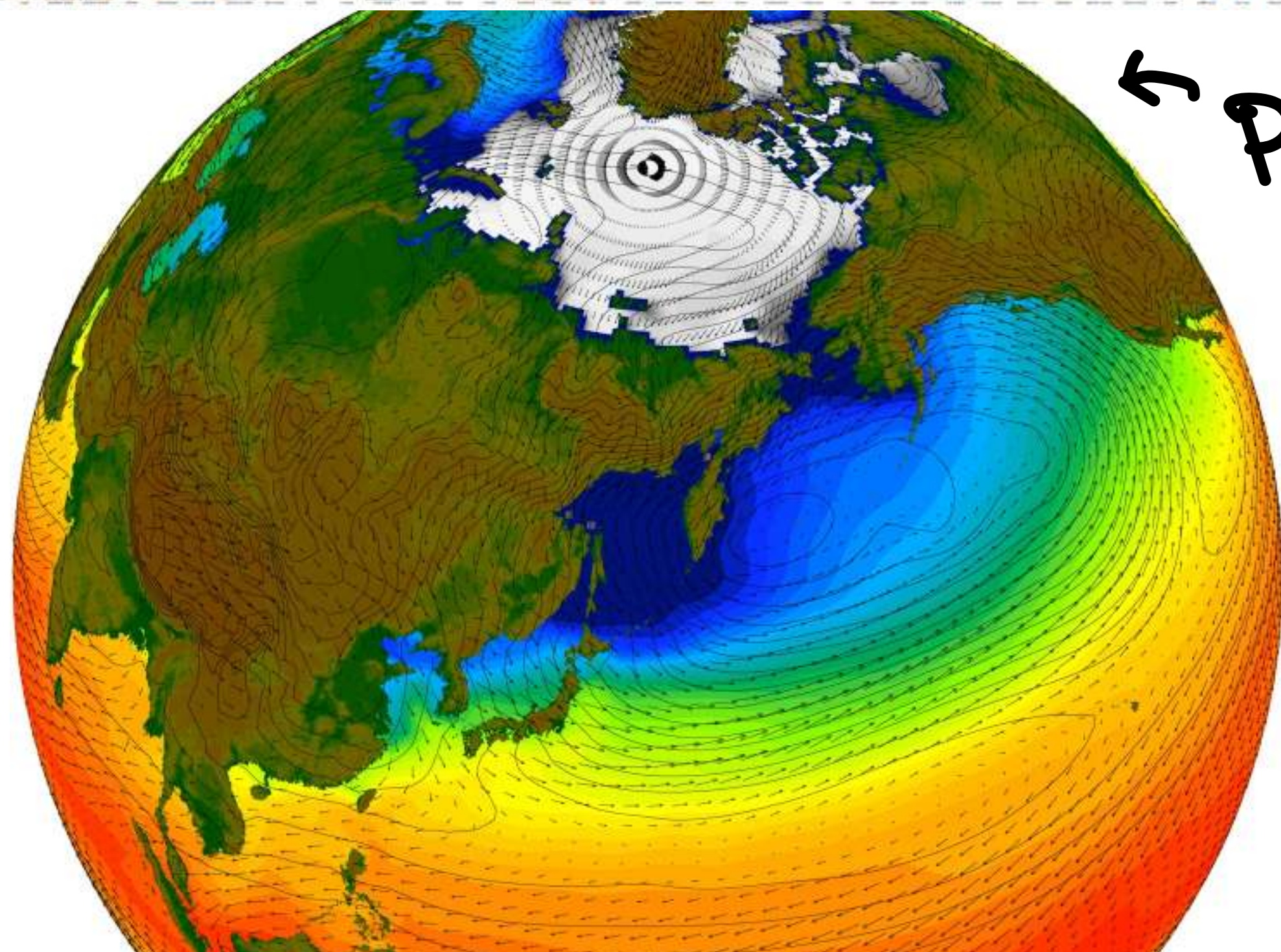
Luis Cerve,
UW



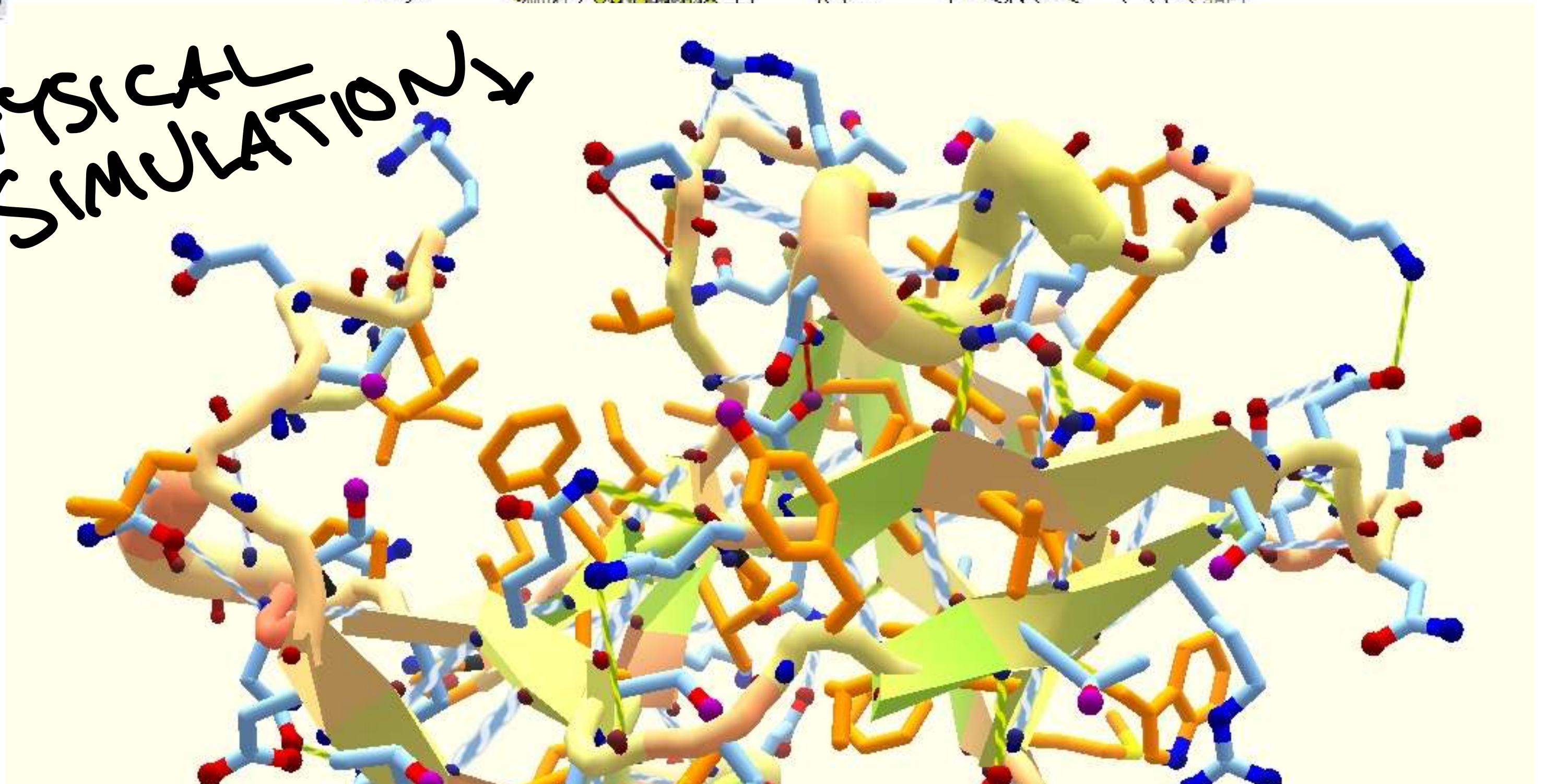
ML

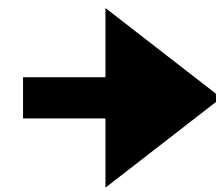
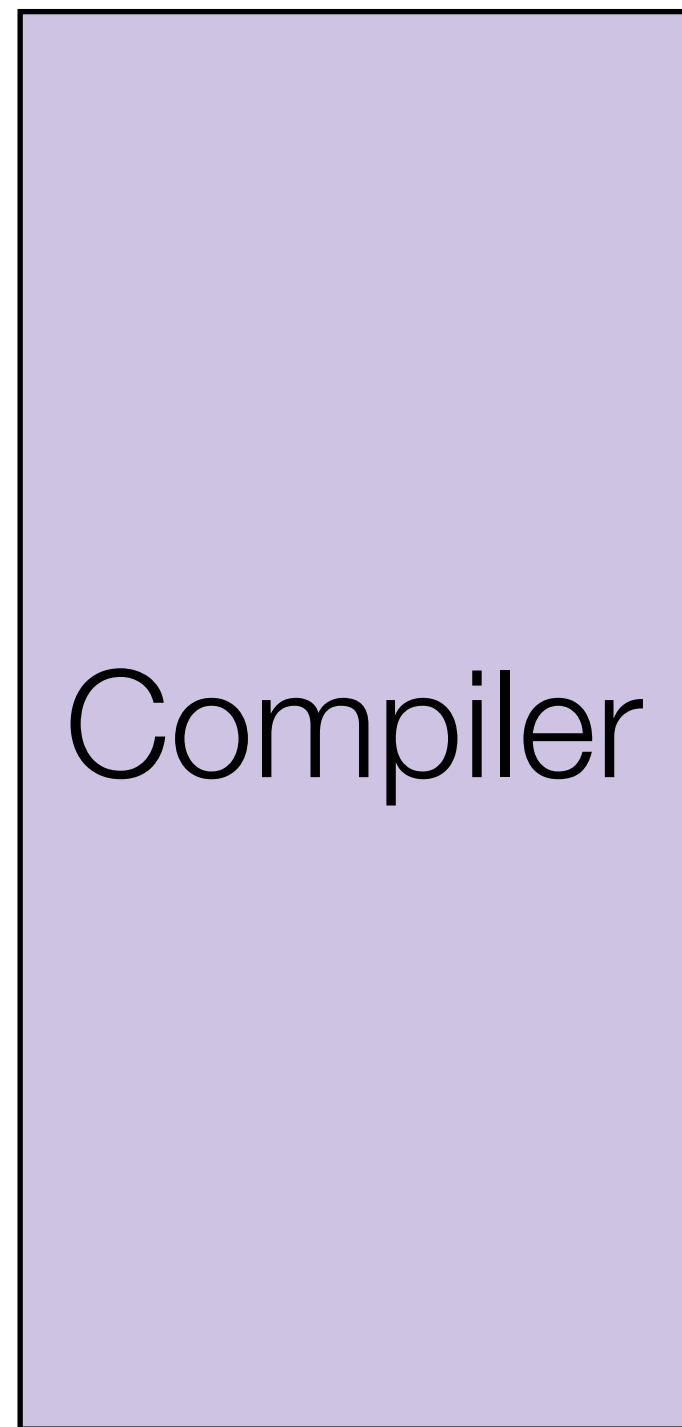


FUZZY SEARCH

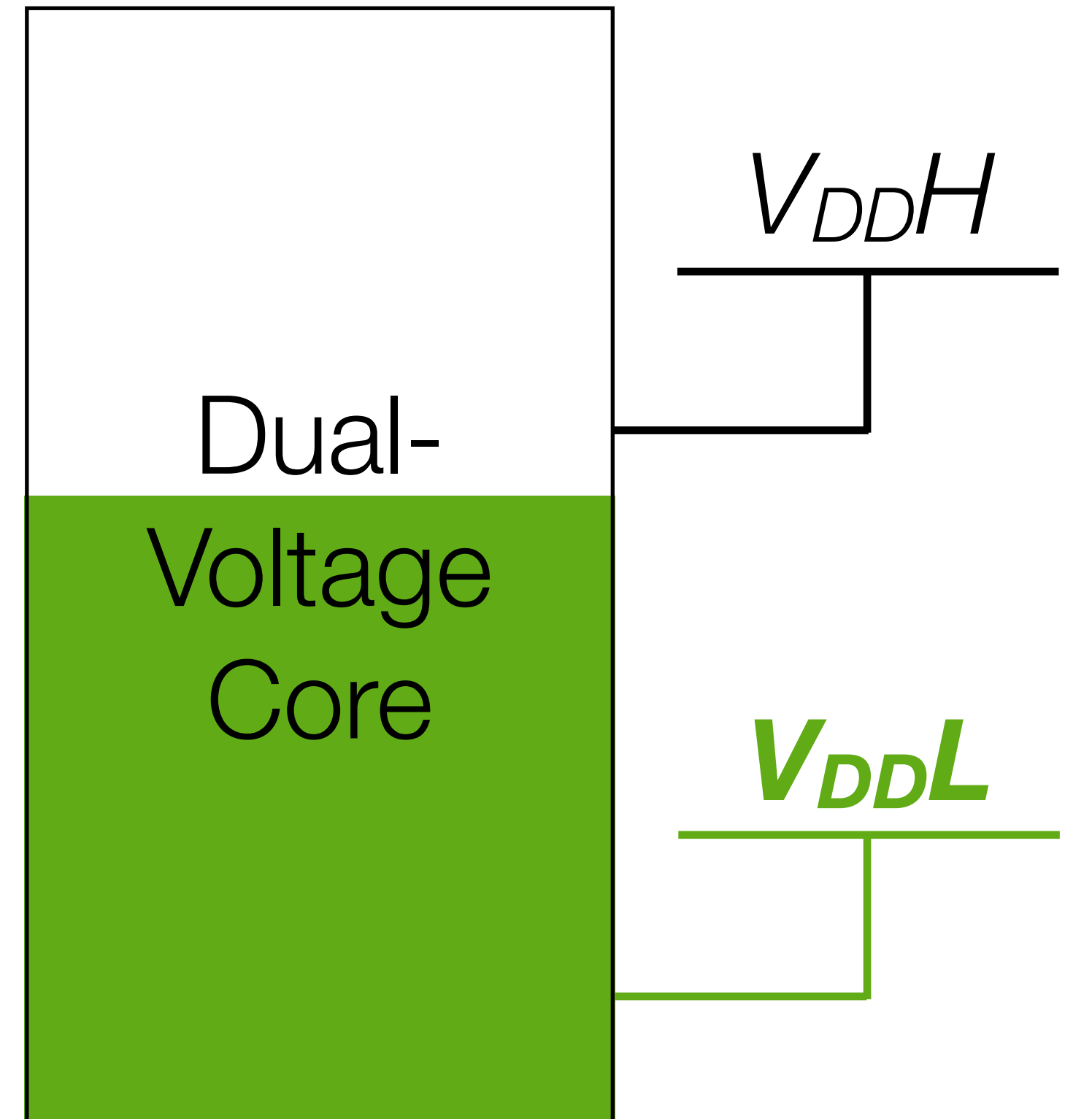
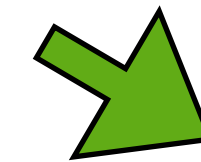
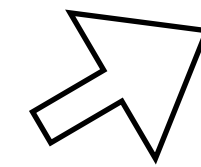


PHYSICAL SIMULATION



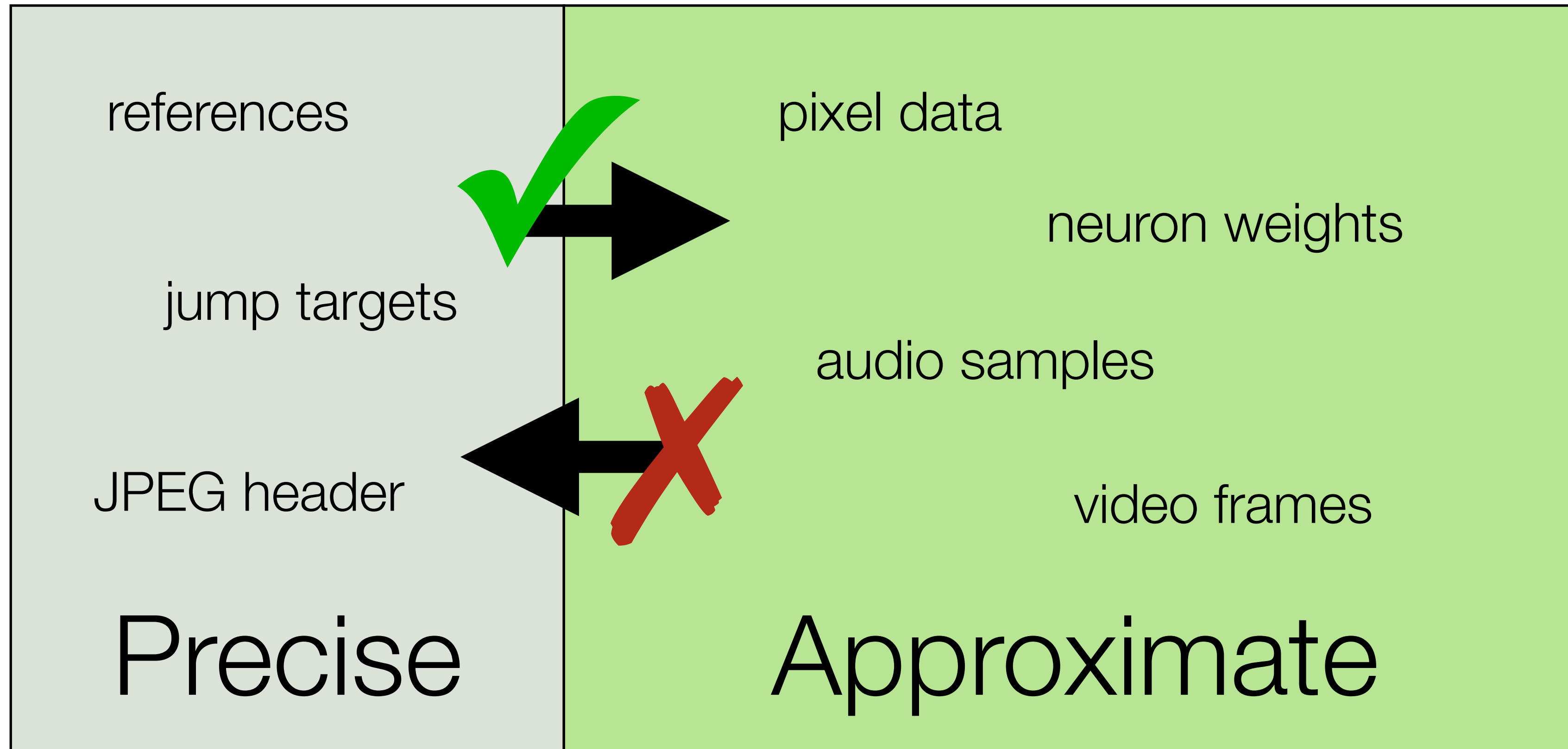


```
ld    0x04  r1
ld    0x08  r2
add.a r1    r2    r3
st.a  0x0c  r3
```



Safety by isolation

[PLDI 2011]





EnerJ type qualifiers

[PLDI 2011]

```
@Approx int a = ...;
```

```
@Precise int p = ...;
```

 `p = a;`


 `a = p;`


EnerJ type qualifiers

[PLDI 2011]

```
@Approx int a = ...;
```

```
@Precise int p = ...;
```

 `p = a;`

 `a = p;`

Operator overloading for approximate computation

```
@Approx int a = ...;
```

```
@Precise int p = ...;
```

```
p + p;
```

precise addition

```
p + a;
```

```
a + a;
```

approximate addition

1. NEVER LISTEN TO ADVICE
2. ALWAYS LISTEN TO ADVICE 初心
3. IT'S NOT NAIVETÉ; IT'S BEGINNER'S MIND
4. THE BEST COMPUTER ARCHITECTURE RESEARCH ALWAYS "CHEATS"
5. DON'T DO FRAMEWORKS
6. START A BLOG
7. TALK TO YOUR FRIENDS

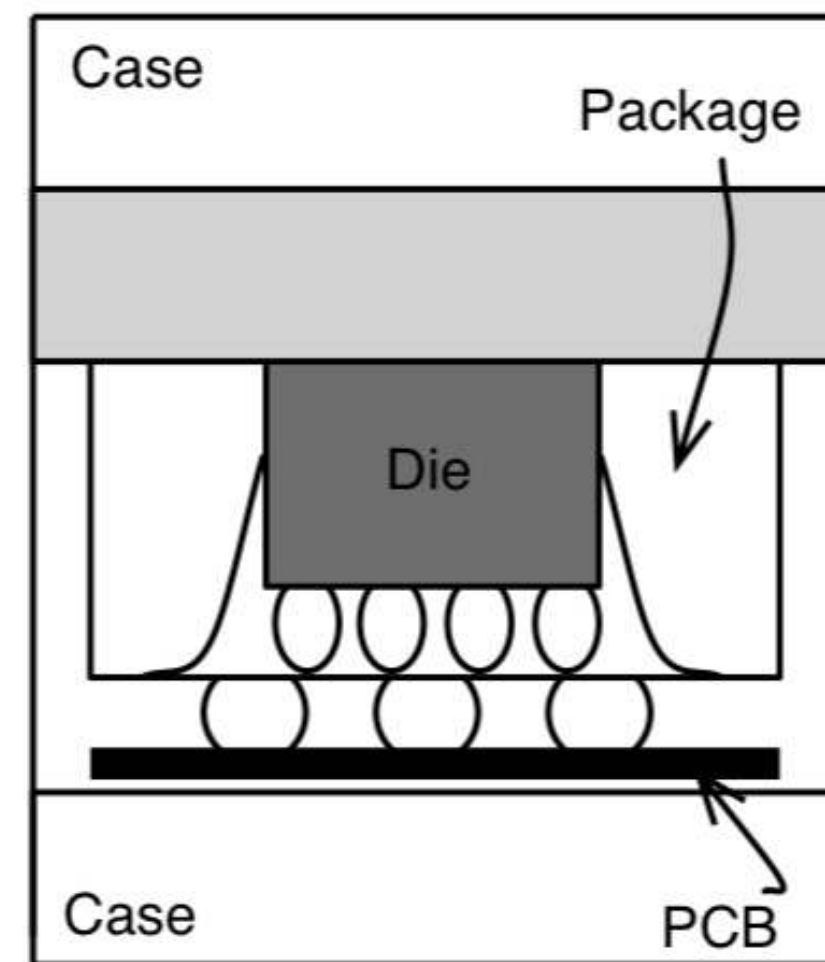
In the Proceedings of the

Marios Papaefthymi

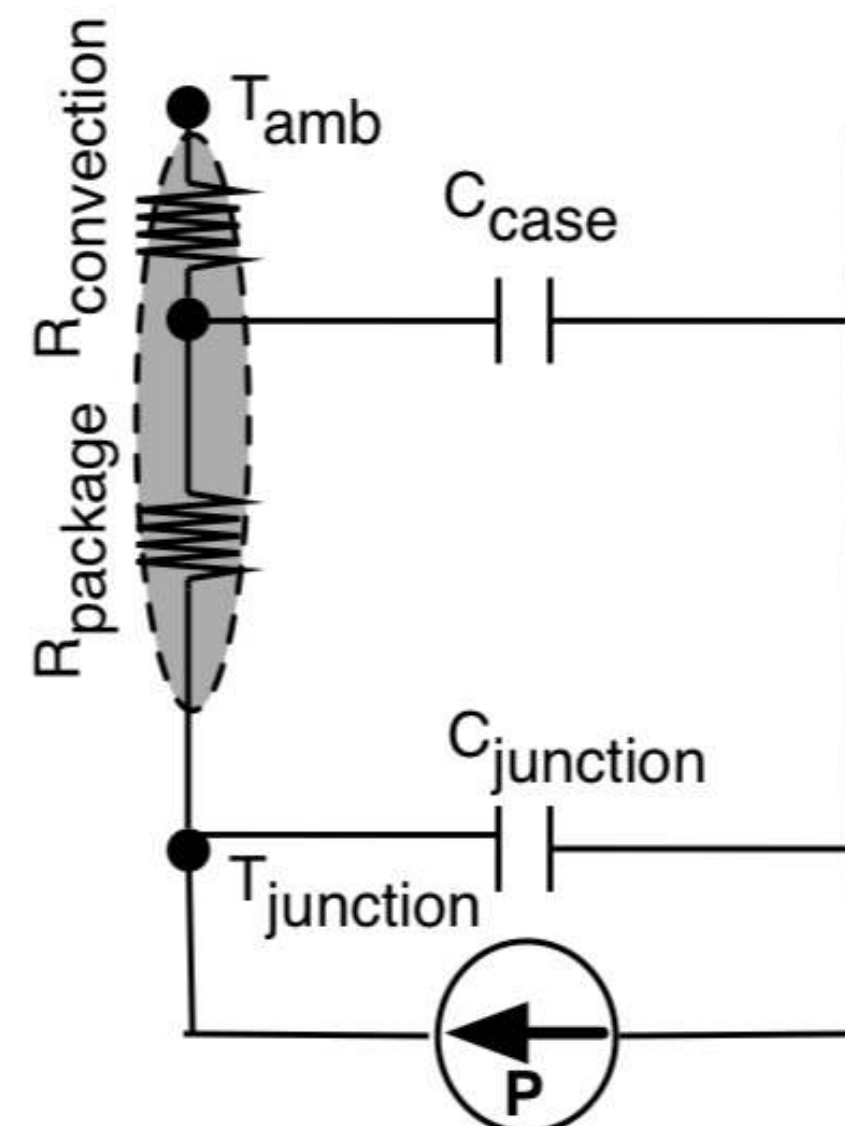
*Department of

†Department of El

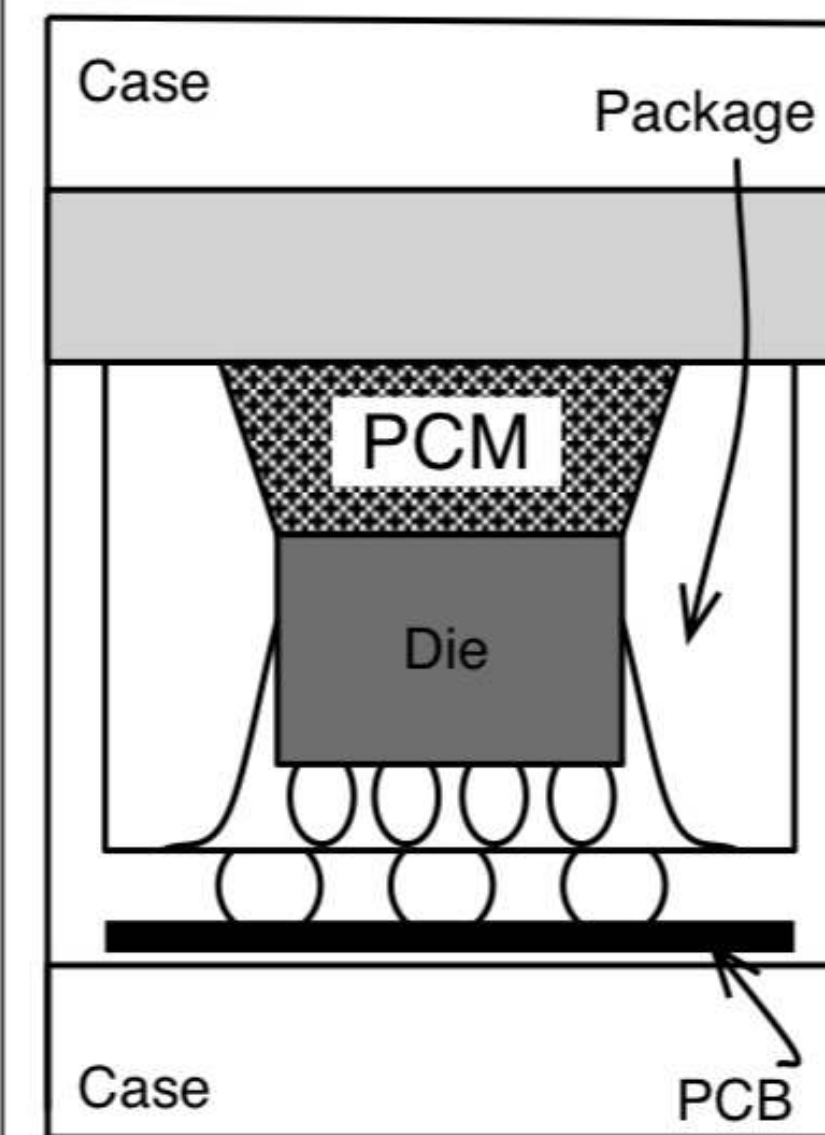
‡Depart



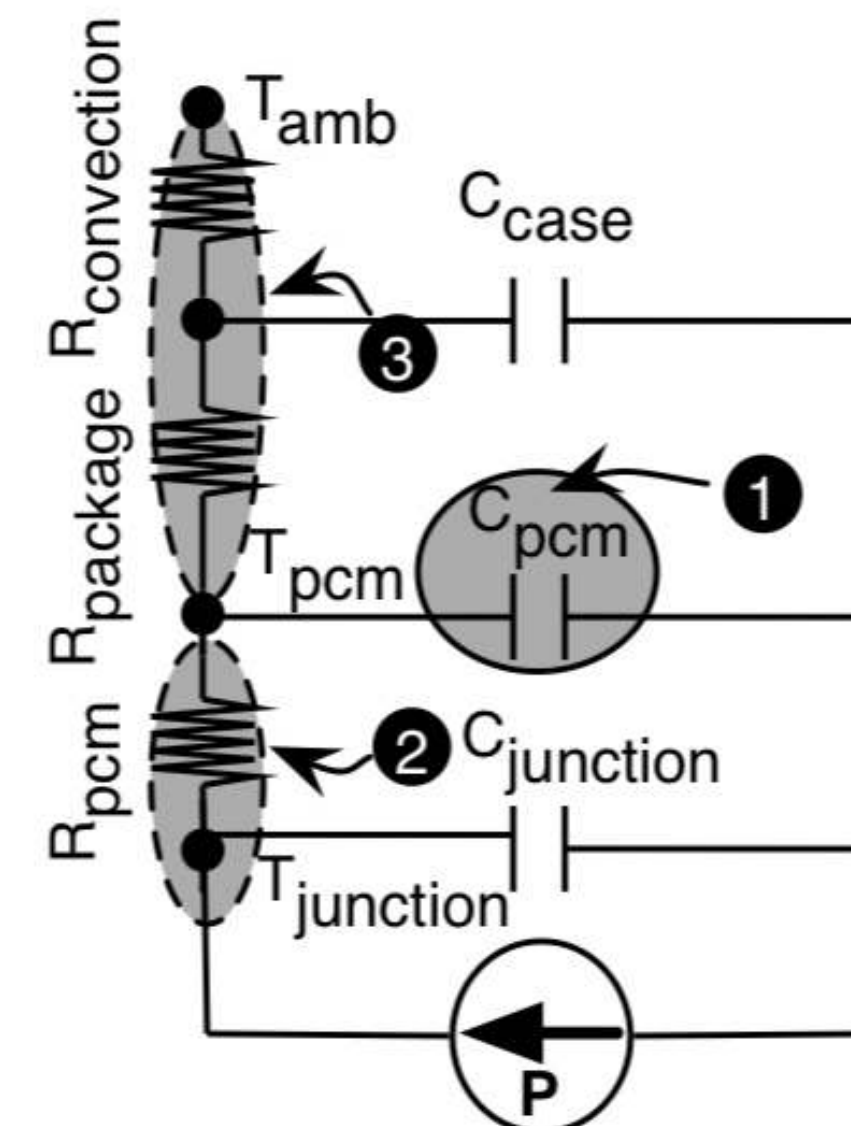
(a)



(b)



(c)



(d)

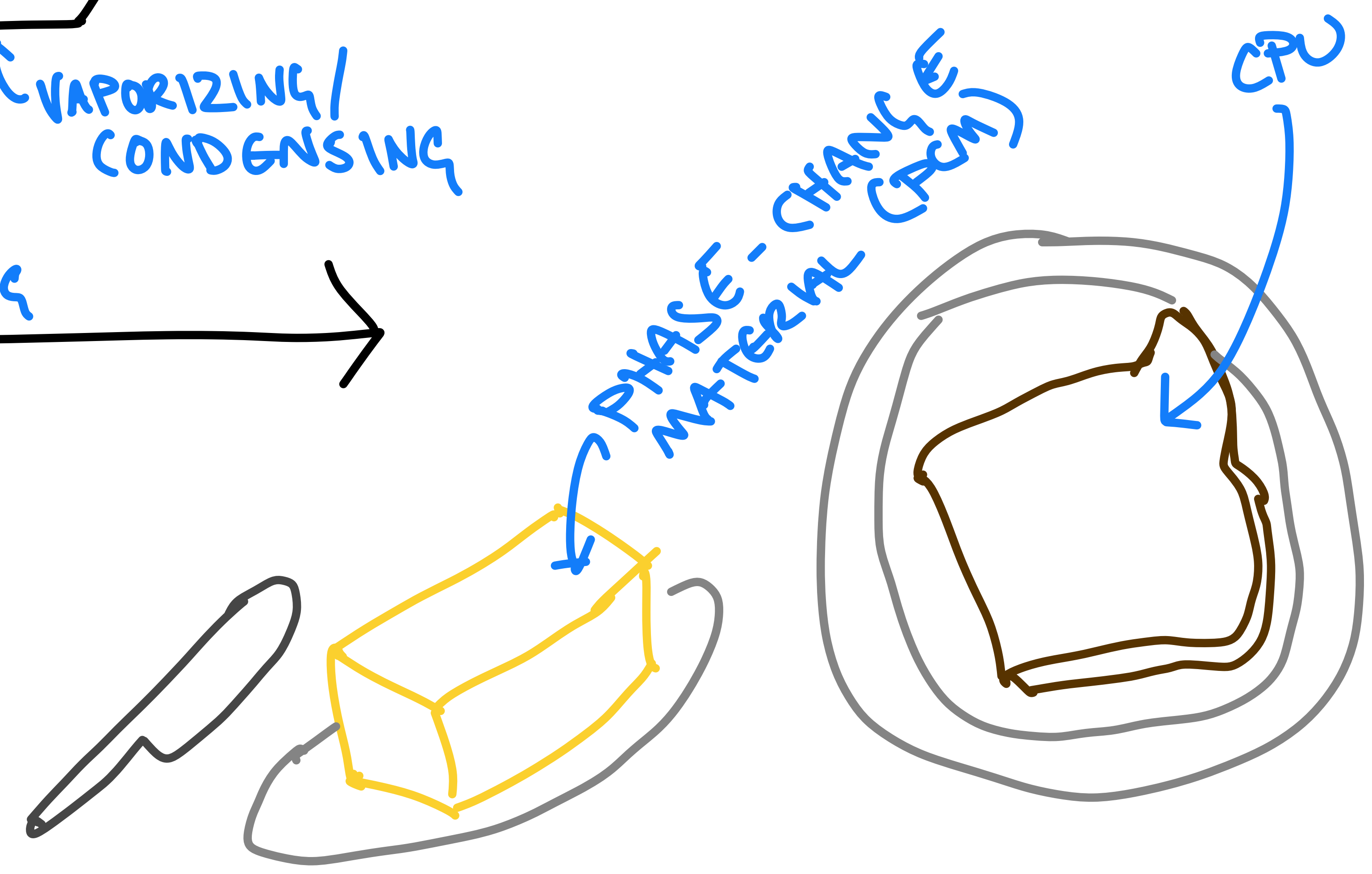
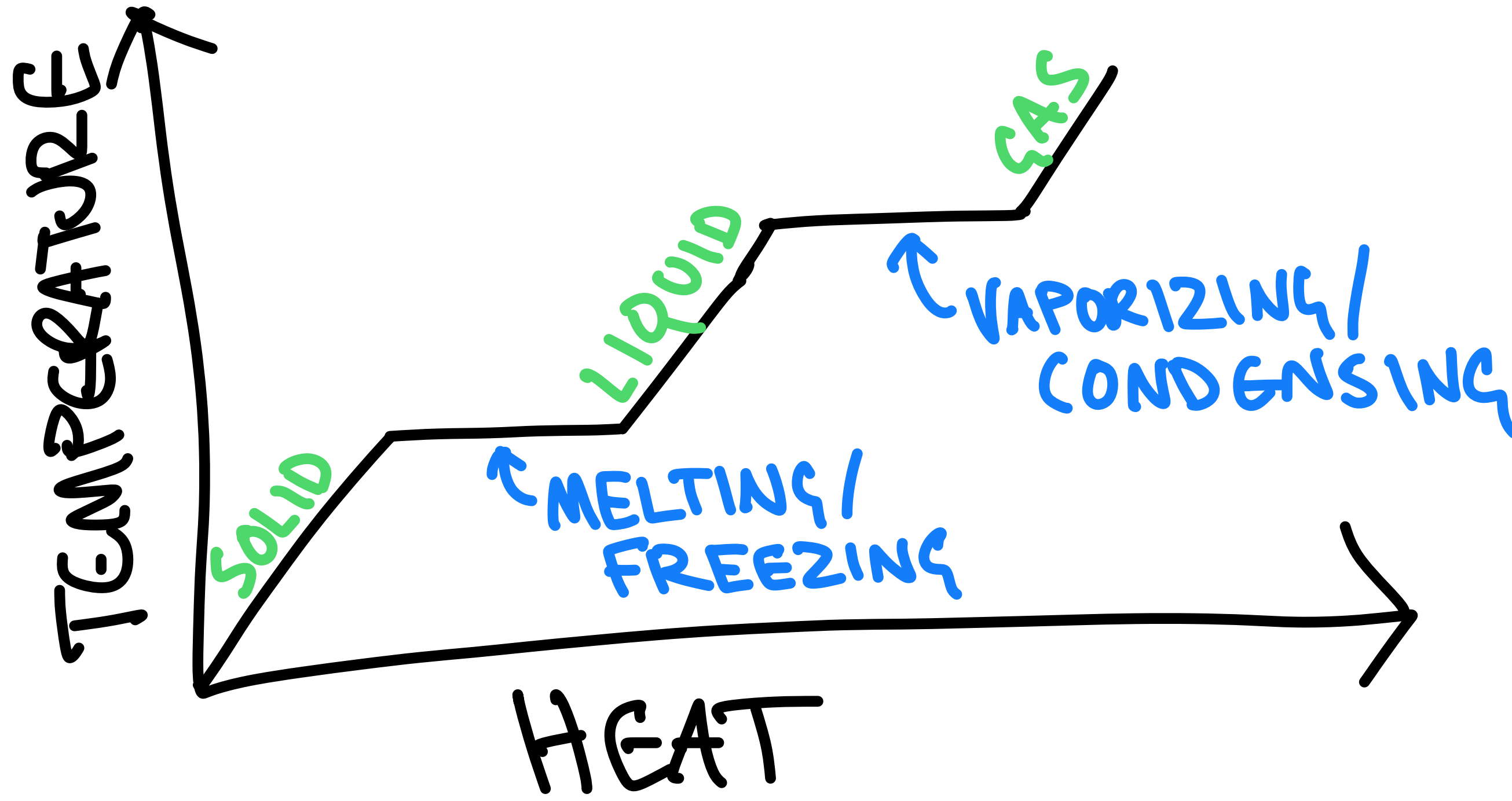
Abstract

Although transistor density continues to increase, voltage scaling has stalled and thus power density is increasing each technology generation. Particularly in mobile devices, which have limited cooling options, these trends lead to a utilization wall in which sustained chip performance is limited primarily by power rather than area. However, many mobile applications do not demand sustained performance; rather they comprise short bursts of computation in response to sporadic user activity.

To improve responsiveness for such applications, this

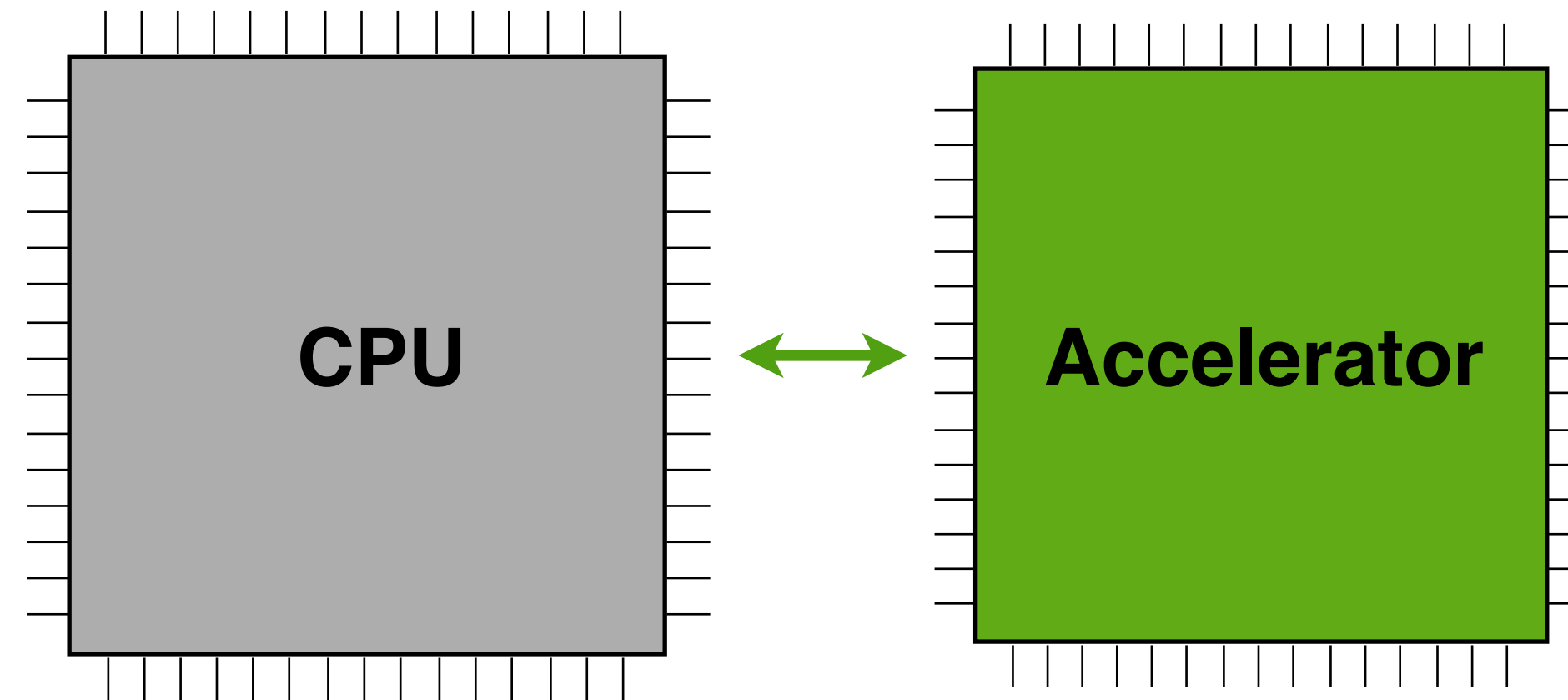
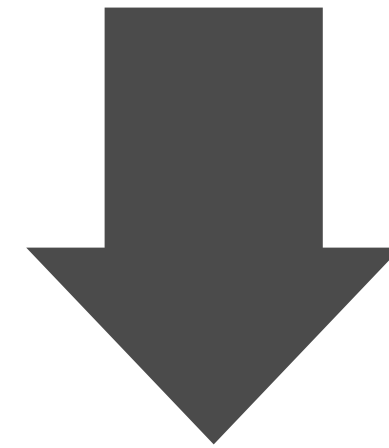
tablets, and smart phones, which, despite relatively large chip area, have thermal budgets that are constrained by the poor heat dissipation of passive convection.

Processors today (and their heat sinks and energy delivery systems) are designed primarily for *sustained performance*. Although the focus on sustained performance is the right design choice for some applications (for example, batch-mode high-performance computing), many workloads are interactive in nature and thus demand *responsiveness*—how long does the user have to wait after initiating a command? In such settings, responsiveness may be more important than sustained



1. NEVER LISTEN TO ADVICE
2. ALWAYS LISTEN TO ADVICE
3. IT'S NOT NAÏVETÉ; IT'S BEGINNER'S MIND
4. THE BEST COMPUTER ARCHITECTURE RESEARCH ALWAYS "CHEATS"
5. DON'T DO FRAMEWORKS
6. START A BLOG
7. TALK TO YOUR FRIENDS

**Approximate
Program**

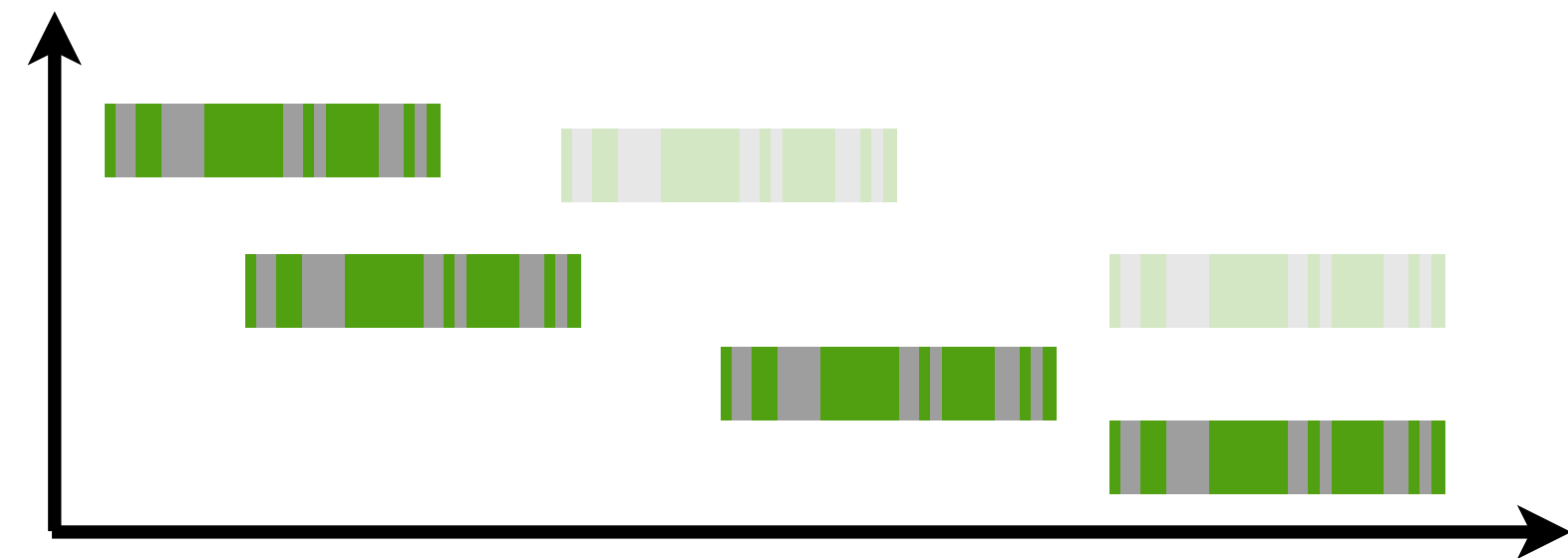


Approximate Program

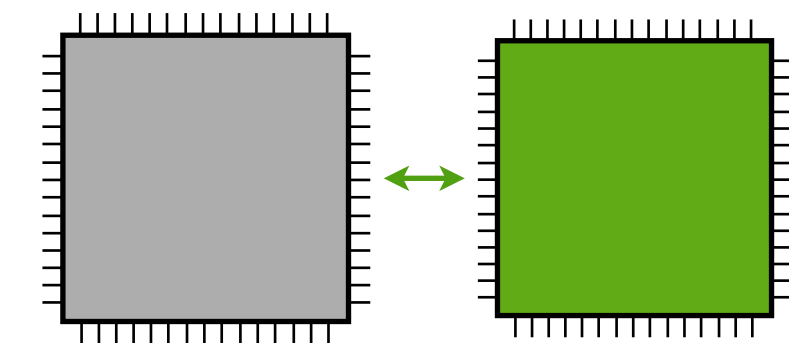
Code with Annotations

+ $d(y, y')$ quality metric

qualifiers for C & C++
interactive feedback



analysis & optimization
auto-tuning



FPGA accelerator

EnerJ type qualifiers for the Clang compiler

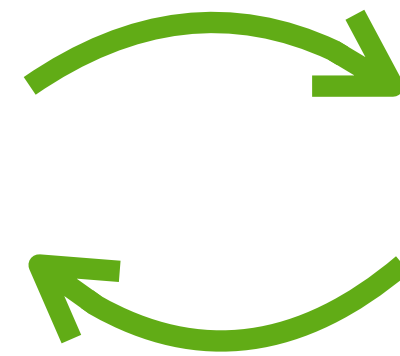
```
for (int k = k1; k < k2; k++) {  
    APPROX float distance =  
        dist(points->p[k],  
            points->p[0]);  
    ...  
}
```


Optimization feedback loop

analysis library finds coarse-grain, safe-to-approximate regions

client optimizations use analysis to relax approximable code

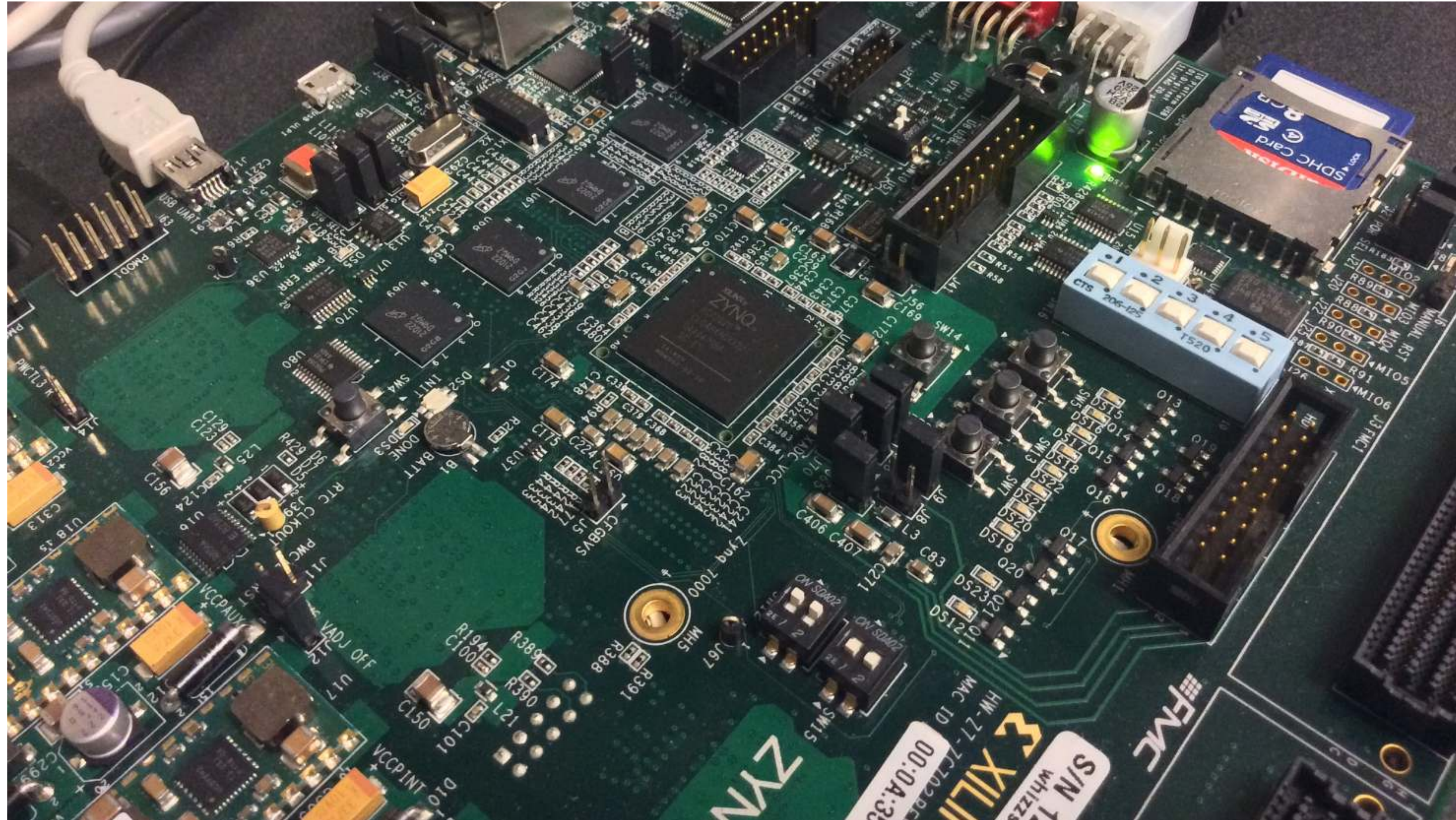
```
for (int k = k1; k < k2; k++) {  
  APPROX float distance =  
    dist(points->p[k],  
         points->p[0]);  
  ...  
}
```



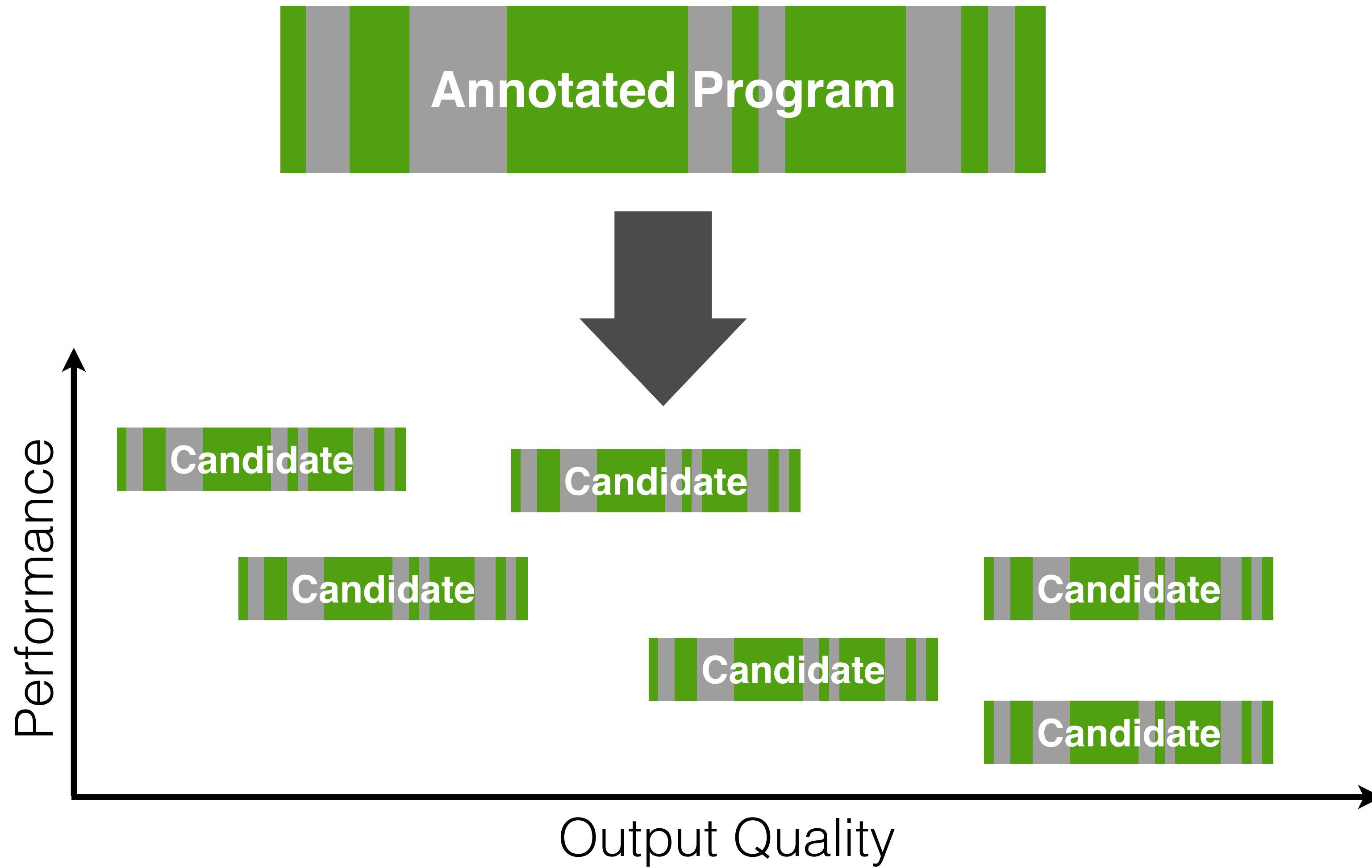
```
$ accept log  
loop at streamcluster.cpp:651  
can offload to NPU  
$
```

reports tell developers where annotations are preventing optimization

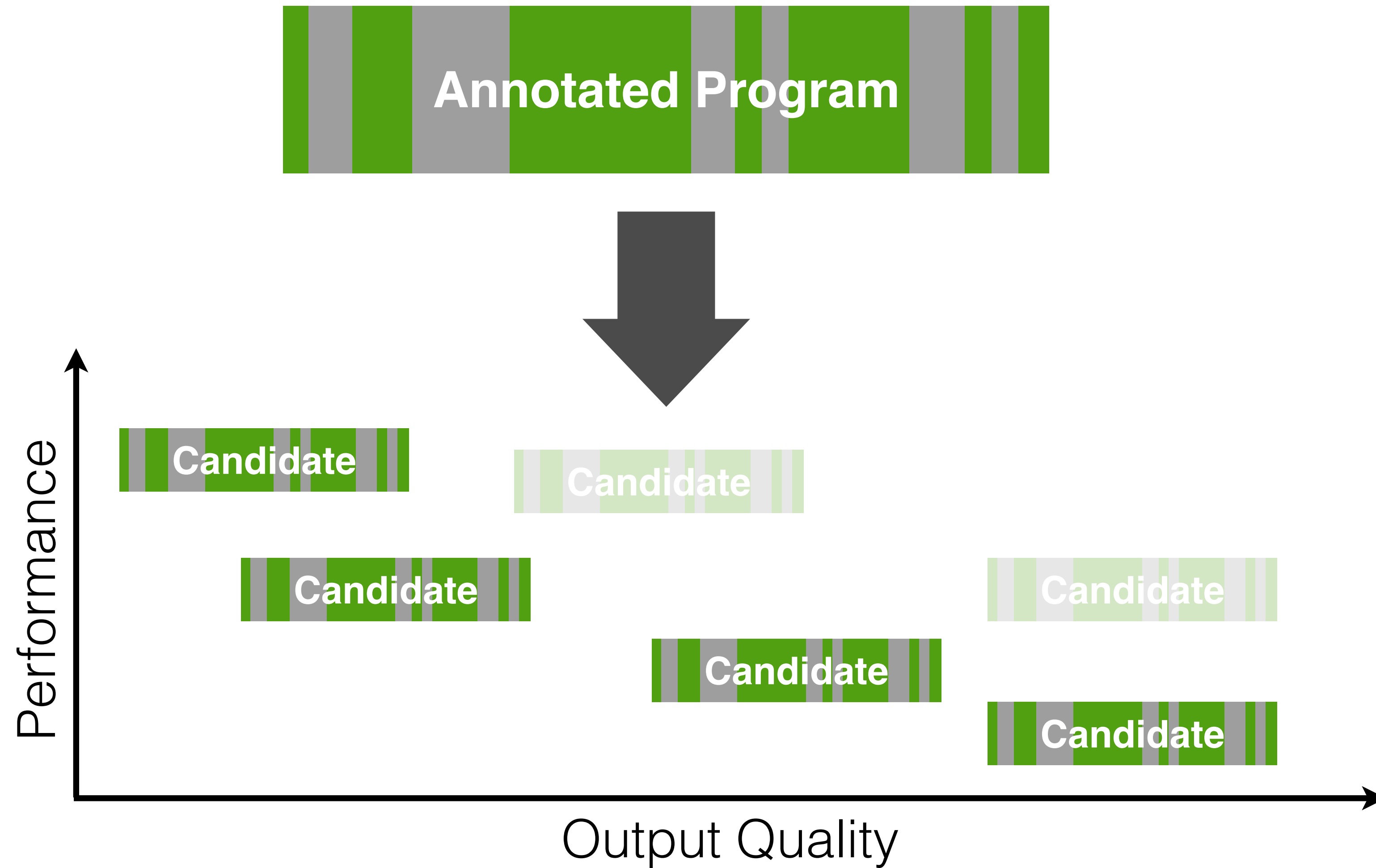
Neural acceleration on a commercial FPGA



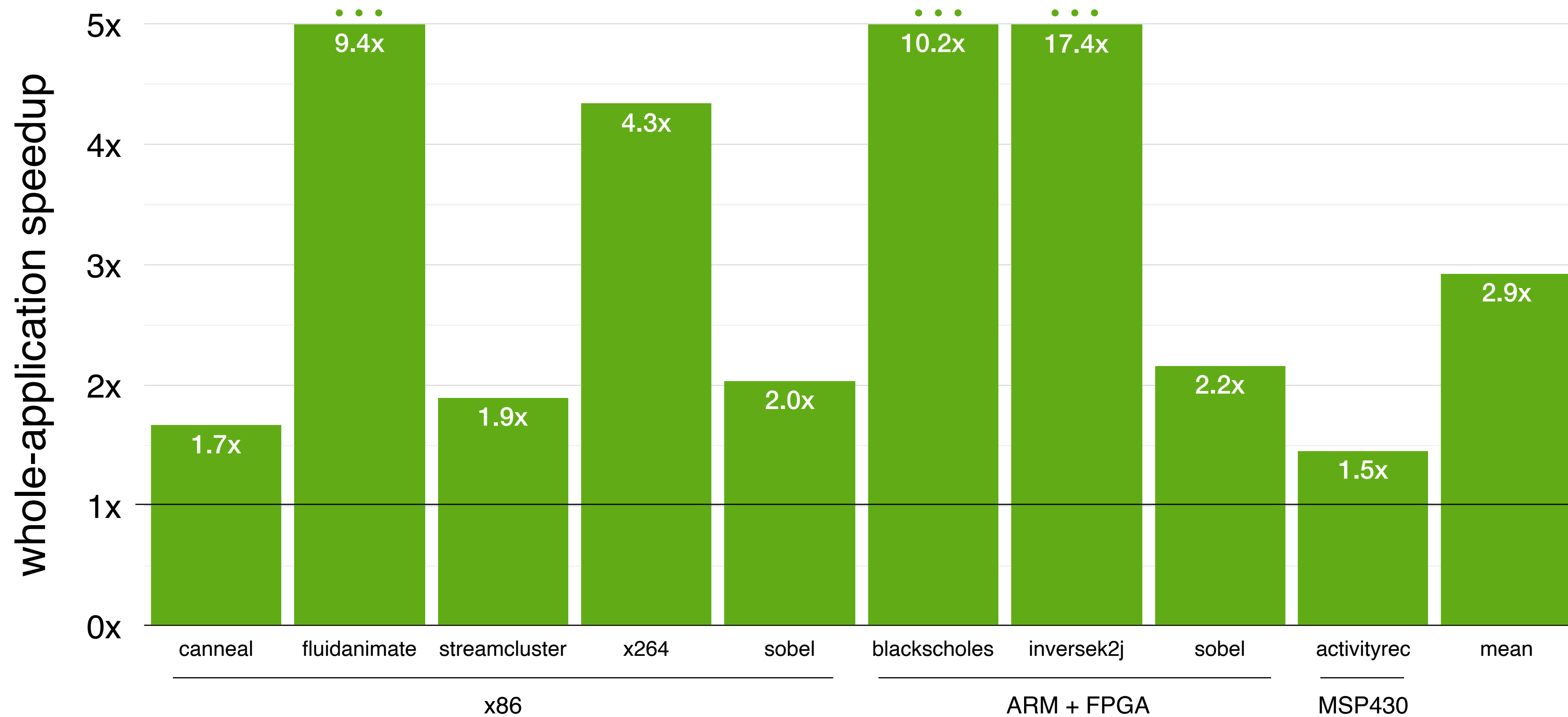
Auto-tuning for optimal trade-offs



Auto-tuning for optimal trade-offs



ACCEPT speedup over precise execution



Benchmarks are 2.9x faster on average
with quality loss under 10%

ASPLOS
PLDI
ASPLOS
DPSLA
CGO
ICS

ACCEPT: A Programmer-Guided Compiler Framework for Practical Approximate Computing

Adrian Sampson André Baixo Benjamin Ransford Thierry Moreau
Joshua Yip Luis Ceze Mark Oskin
University of Washington

Abstract

Approximate computing trades off accuracy for better performance and energy efficiency. It offers promising optimization

sion [Rubio-González et al. 2013]; using special low-power hardware structures that produce wrong results probabilistically [Esmailzadeh et al. 2012b; Liu et al. 2011]; and training hardware neural networks to mimic the behavior of costly

1. NEVER LISTEN TO ADVICE
2. ALWAYS LISTEN TO ADVICE
3. IT'S NOT NAÏVETÉ; IT'S BEGINNER'S MIND
4. THE BEST COMPUTER ARCHITECTURE RESEARCH ALWAYS "CHEATS"
5. DON'T DO FRAMEWORKS
6. START A BLOG
7. TALK TO YOUR FRIENDS

Closed Problems in Approximate Computing

OCTOBER 14, 2017

These are notes for a talk I will give at the NOPE workshop at MICRO 2017, where the title is *Approximate Computing Is Dead; Long Live Approximate Computing*.

Approximate computing has reached an adolescent phase as a research area. We have picked bushels of low-hanging fruit. While there are many approximation papers left to write, it's a good time to enumerate the closed problems: research problems that are probably no longer worth pursuing.

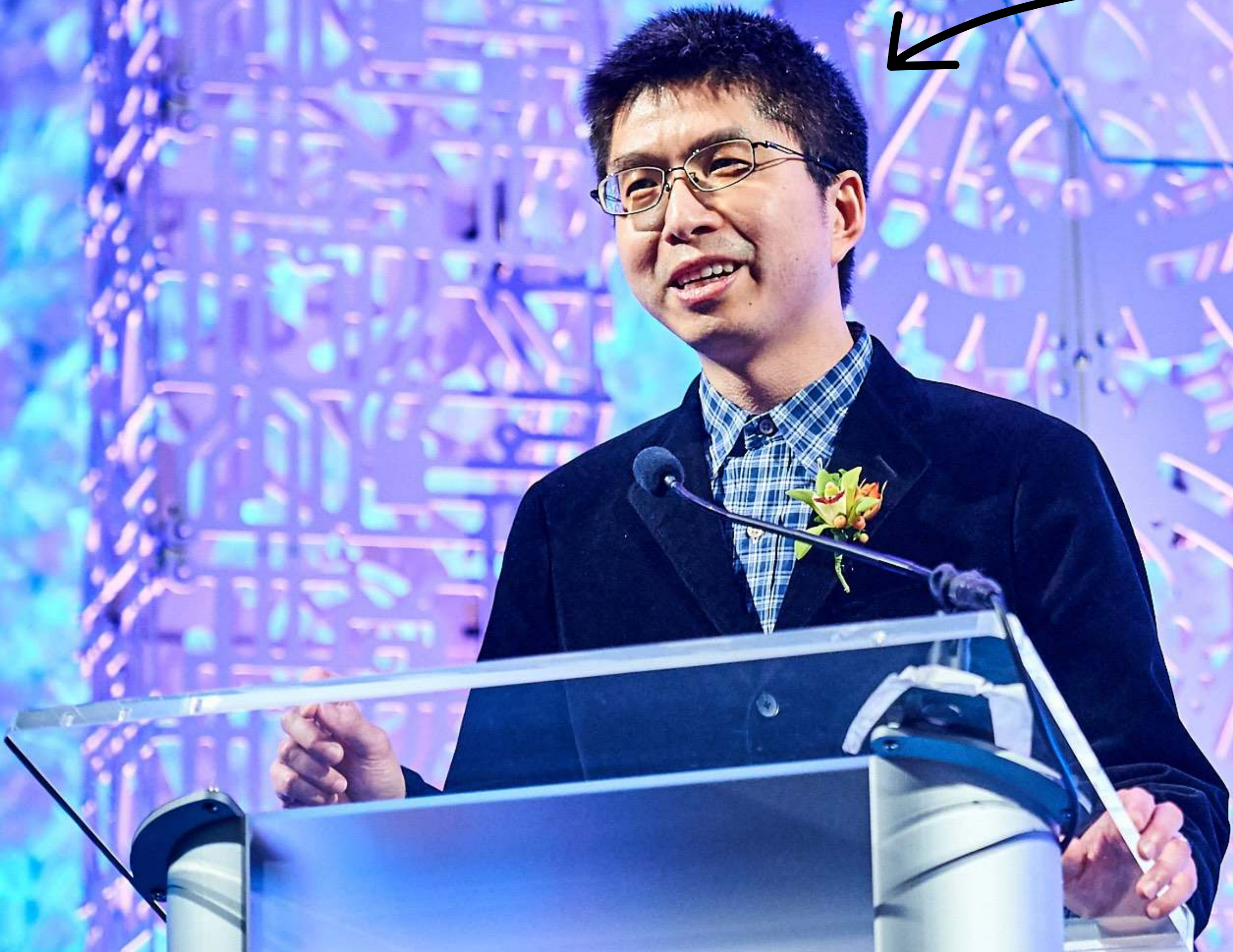
Closed Problems in Approximate Hardware

No more approximate functional units. Especially for people who love VLSI work, a natural first step in approximate computing is to design approximate adders, multipliers, and other basic functional units. Cut a carry chain here, drop a block of intermediate results there, or use an automated search to find “unnecessary” gates—there are lots of ways to design an FU that's mostly right most of the time. Despite dozens of papers in this

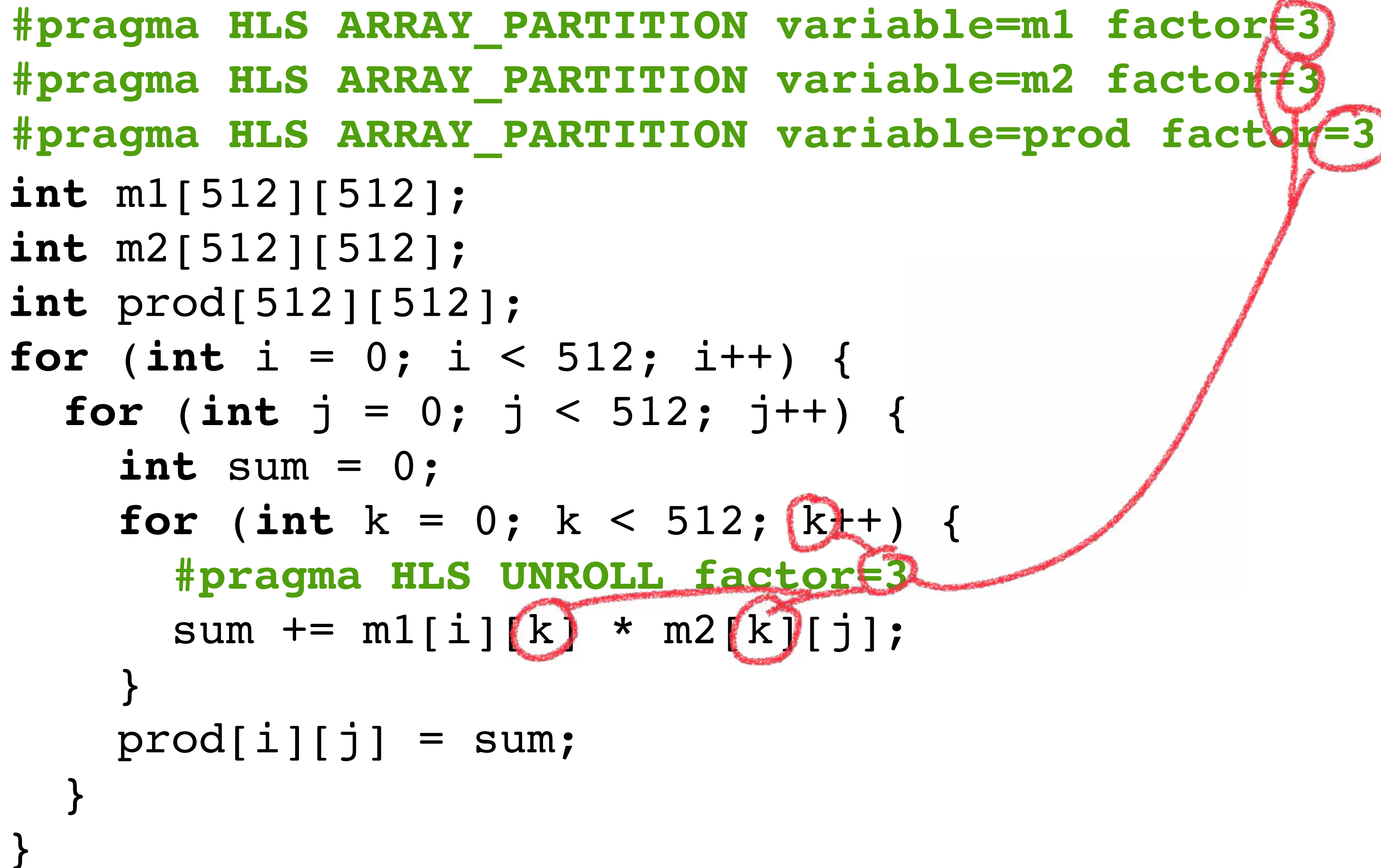
- PRACTICE WRITING
- CLARIFY YOUR OWN THOUGHTS
- SAVE TIME EXPLAINING STUFF
- “MARKET” YOUR WORK
- GET FAMOUS

1. NEVER LISTEN TO ADVICE
2. ALWAYS LISTEN TO ADVICE
3. IT'S NOT NAÏVETÉ; IT'S BEGINNER'S MIND
4. THE BEST COMPUTER ARCHITECTURE RESEARCH ALWAYS "CHEATS"
5. DON'T DO FRAMEWORKS
6. START A BLOG
7. TALK TO YOUR FRIENDS

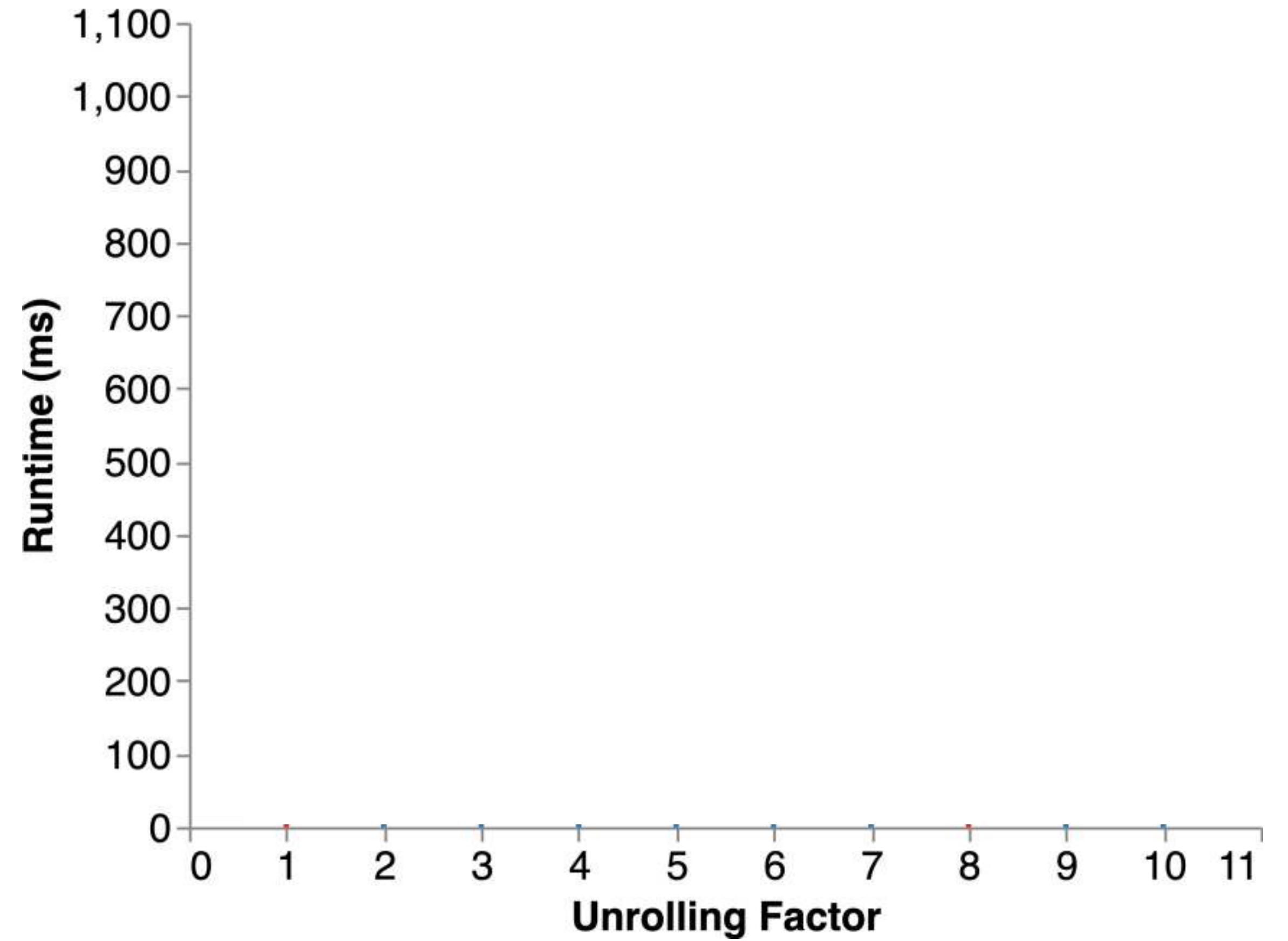
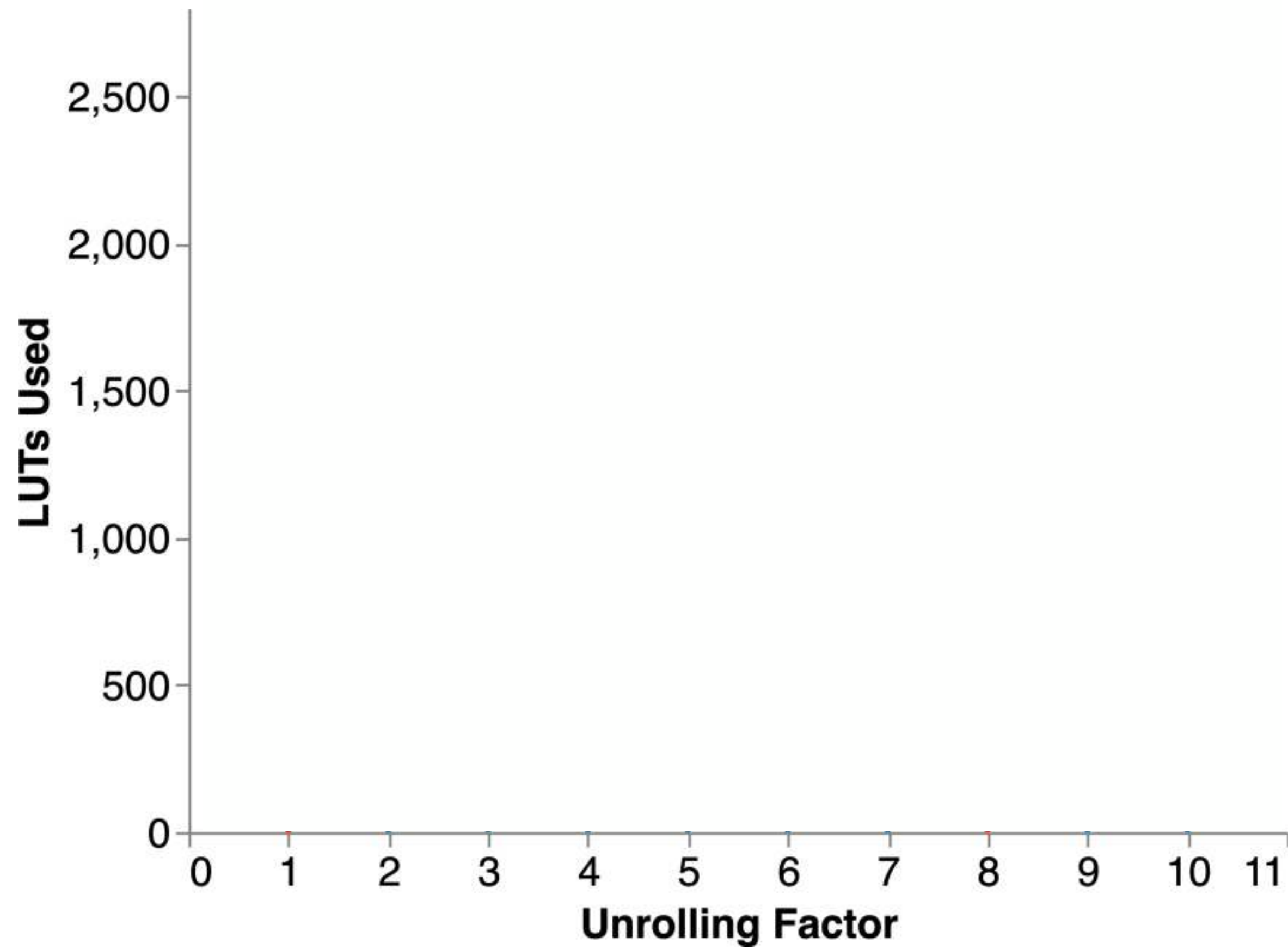
ZHRU ZHANG
CORNELL



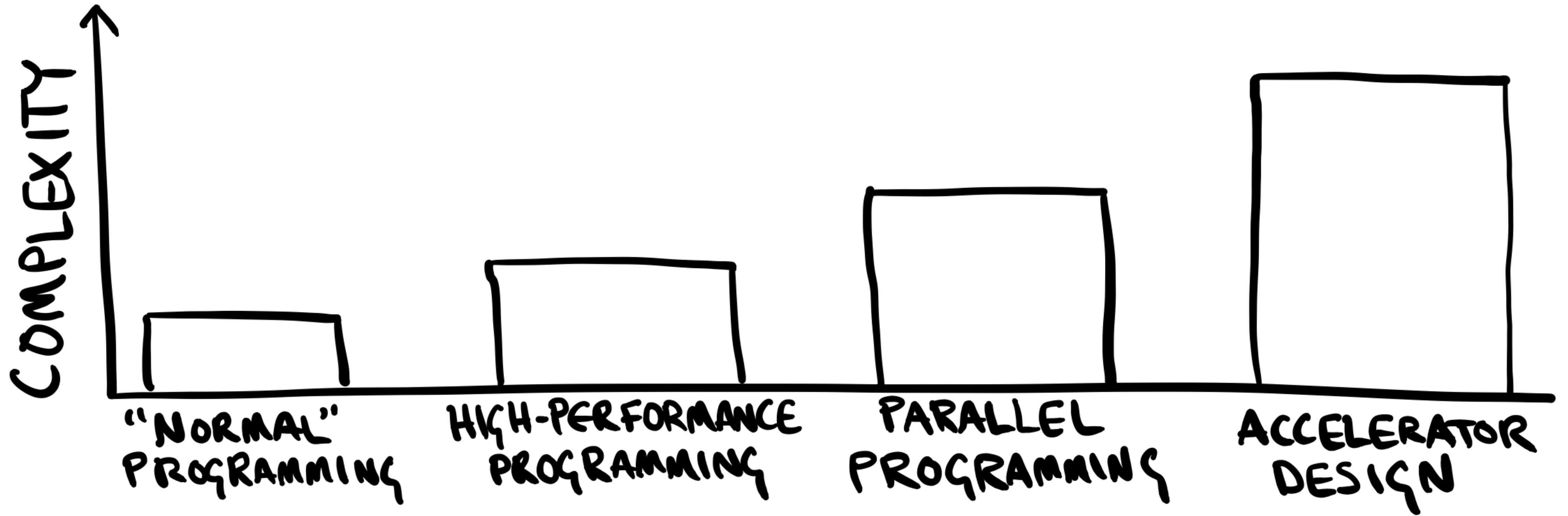

```
#pragma HLS ARRAY_PARTITION variable=m1 factor=3
#pragma HLS ARRAY_PARTITION variable=m2 factor=3
#pragma HLS ARRAY_PARTITION variable=prod factor=3
int m1[512][512];
int m2[512][512];
int prod[512][512];
for (int i = 0; i < 512; i++) {
    for (int j = 0; j < 512; j++) {
        int sum = 0;
        for (int k = 0; k < 512; k++) {
            #pragma HLS UNROLL factor=3
            sum += m1[i][k] * m2[k][j];
        }
        prod[i][j] = sum;
    }
}
```

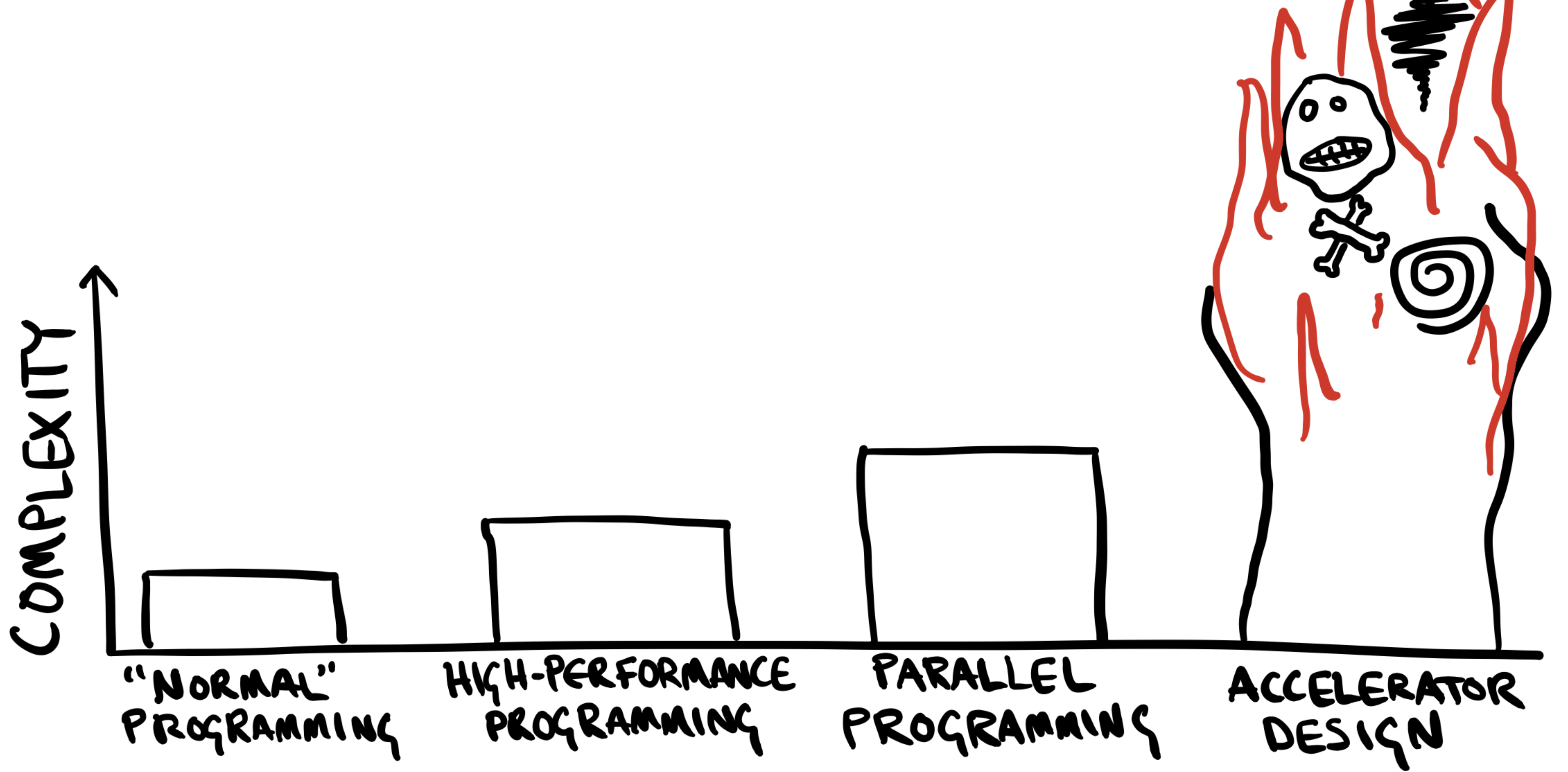


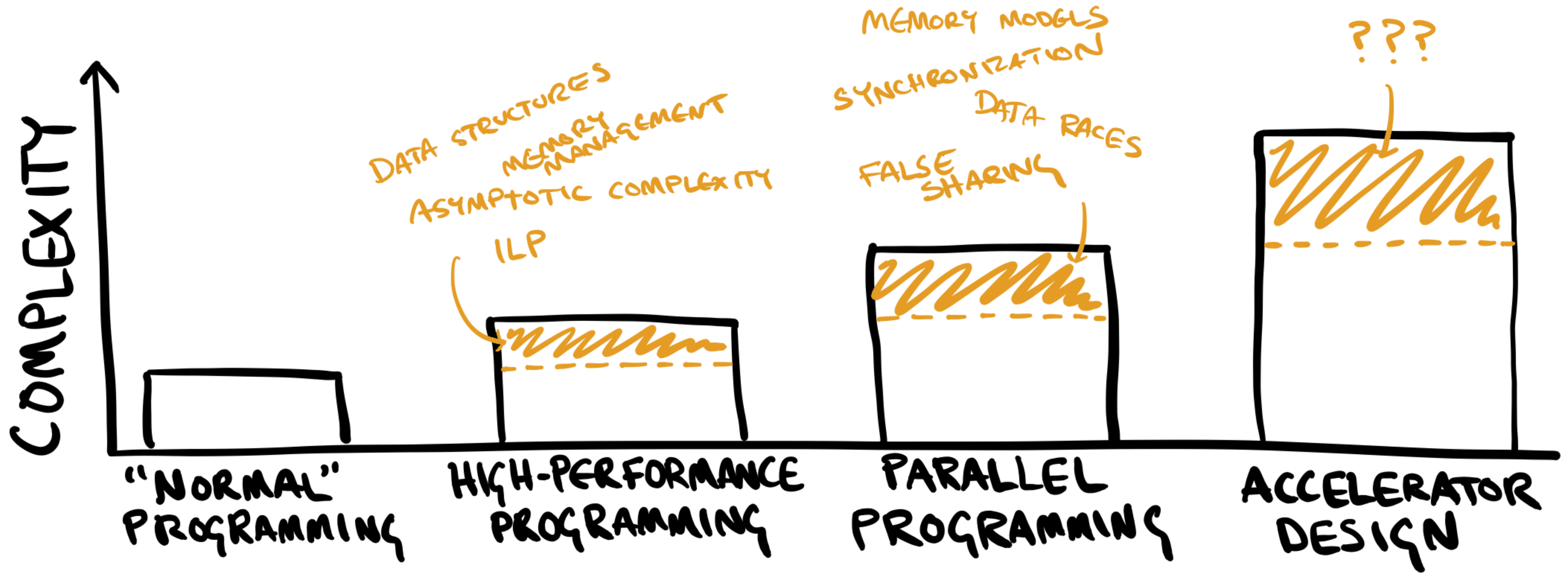
Area-Performance Trade-Offs in Unrolling



1. NEVER LISTEN TO ADVICE
2. ALWAYS LISTEN TO ADVICE
3. IT'S NOT NAÏVETÉ; IT'S BEGINNER'S MIND
4. THE BEST COMPUTER ARCHITECTURE RESEARCH ALWAYS "CHEATS"
5. DON'T DO FRAMEWORKS
6. START A BLOG
7. TALK TO YOUR FRIENDS








```

layer_kernels.cu + (~/Downloads/cuda-convnet/trunk/src) - VIM3
/*
 * E = -log(y_t)
 * y_l: (numOut, numCases)
 * labels: (1, numCases)
 *
 * dE_dx_l: (numOut, numCases)
 */
template <bool ad...
__global__ void k...__global__ void kDotProduct_r...
int numCases, ~ols, const uint numElements)
__shared__ float shmem[DP...
const int tx...
const int ty...
const int tid...
uint eid = DP_BLOCKSIZE +...
shmem[threadIdx.x] = 0;
if (eid < numCol...
for (; eid <...

gpumodel.py + (~/Downloads/cuda-convnet/trunk) - VIM9
# GPU Model interface
class IGPUModel:
    def __init__(self, model_name, op, load_dic, filename_options=None, dp_par...
        ~ams={}):
        # these are input parameters
        self.model_name = model_name
        self.op = op
        self.options = op.options
        self.load_dic = load_dic
        self.filename_options = filename_options

matrix.cpp + (~/Downloads/cuda-convnet/trunk/src/common) - VIM8
void Matrix::updateDims(long int numRows, long int numCols) {
    this->_numRows = numRows;
    this->_numCols = numCols;
    this->_numElements = numRows * numCols;
}

void Matrix::checkBounds(long int startRow, long int endRow, long int startCo...
    ~l, long int endCol) const {
    assert(startRow >= 0 && startRow <= _numRows);
    assert(endRow >= 0 && endRow <= _numRows);
    assert(startCol >= 0 && startCol <= _numCols);
    assert(endCol >= 0 && endCol <= _numCols);
}

/* will return a view if possible */
Matrix& Matrix::slice(long int startRow, long int endRow, long int startCol, l...
    ~ong int endCol) const {
    endRow = endRow < 0 ? this->_numRows : endRow;
    endCol = endCol < 0 ? this->_numCols : endCol;
    _checkBounds(startRow, endRow, startCol, endCol);
    if (!isTrans() && ((startCol == 0 && endCol == this->_numCols) || (startRo...
        ~w == endRow - 1))) {
        return *new Matrix(this->_data + startRow * this->_numCols + startCol,
            ~ endRow - startRow, endCol - startCol);
    } else if (isTrans() && ((startRow == 0 && endRow == this->_numRows) || (s...
        ~tartCol == endCol - 1))) {
        return *new Matrix(this->_data + startCol * this->_numRows + startRow,
            ~ endRow - startRow, endCol - startCol, true);
    }
    Matrix& newSlice = *new Matrix(endRow - startRow, endCol - startCol);
    this->copy(newSlice, startRow, endRow, startCol, endCol, 0, 0);
    return newSlice;
}

convdata.py + (~/Downloads/cuda-convnet/trunk/src/common) - VIM7
class CroppedCIFARDataProvider:
    def __init__(self, data_dir, labels_dir, num_epochs, num_batches_per_epoch, test...
        ~=None, dp_params=None, test_data_dir=None, labels_dir=None, num_epochs=10, num_batches_per_epoch=10,
        ~poch, init_batchnum, dp_params=None):
        self.border_size = dp_params.border_size
        self.inner_size = 32
        self.multiview = dp_params.multiview
        self.num_views = 5*2
        self.data_mult = self.border_size + self.inner_size
        self.num_colors = 3

        for d in self.data_dirs:
            d['data'] = n.recursive_glob(d, ['data'])
            d['labels'] = n.recursive_glob(d, ['labels'])

        self.cropped_data = self._load_data(self.data_dir, self.data_mult)
        self.cropped_labels = self._load_labels(self.labels_dir, self.data_mult)

        self.batches_generate = self._generate_batches(self.cropped_data, self.cropped_labels,
            ~f.border_size: self.border_size, self.inner_size)

        def get_next_batch(self):
            epoch, batchnum, data, labels = self._get_next_batch()

            cropped = self.cropped_data[epoch][batchnum]
            self.__trim_borders(cropped)
            cropped -= self.data_mult

matrix.cpp + (~/Downloads/cuda-convnet/trunk/src/common) - VIM8
void Matrix::updateDims(long int numRows, long int numCols) {
    this->_numRows = numRows;
    this->_numCols = numCols;
    this->_numElements = numRows * numCols;
}

void Matrix::checkBounds(long int startRow, long int endRow, long int startCo...
    ~l, long int endCol) const {
    assert(startRow >= 0 && startRow <= _numRows);
    assert(endRow >= 0 && endRow <= _numRows);
    assert(startCol >= 0 && startCol <= _numCols);
    assert(endCol >= 0 && endCol <= _numCols);
}

/* will return a view if possible */
Matrix& Matrix::slice(long int startRow, long int endRow, long int startCol, l...
    ~ong int endCol) const {
    endRow = endRow < 0 ? this->_numRows : endRow;
    endCol = endCol < 0 ? this->_numCols : endCol;
    _checkBounds(startRow, endRow, startCol, endCol);
    if (!isTrans() && ((startCol == 0 && endCol == this->_numCols) || (startRo...
        ~w == endRow - 1))) {
        return *new Matrix(this->_data + startRow * this->_numCols + startCol,
            ~ endRow - startRow, endCol - startCol);
    } else if (isTrans() && ((startRow == 0 && endRow == this->_numRows) || (s...
        ~tartCol == endCol - 1))) {
        return *new Matrix(this->_data + startCol * this->_numRows + startRow,
            ~ endRow - startRow, endCol - startCol, true);
    }
    Matrix& newSlice = *new Matrix(endRow - startRow, endCol - startCol);
    this->copy(newSlice, startRow, endRow, startCol, endCol, 0, 0);
    return newSlice;
}

convdata.py + (~/Downloads/cuda-convnet/trunk/src/common) - VIM7
class CroppedCIFARDataProvider:
    def __init__(self, data_dir, labels_dir, num_epochs, num_batches_per_epoch, test...
        ~=None, dp_params=None, test_data_dir=None, labels_dir=None, num_epochs=10, num_batches_per_epoch=10,
        ~poch, init_batchnum, dp_params=None):
        self.border_size = dp_params.border_size
        self.inner_size = 32
        self.multiview = dp_params.multiview
        self.num_views = 5*2
        self.data_mult = self.border_size + self.inner_size
        self.num_colors = 3

        for d in self.data_dirs:
            d['data'] = n.recursive_glob(d, ['data'])
            d['labels'] = n.recursive_glob(d, ['labels'])

        self.cropped_data = self._load_data(self.data_dir, self.data_mult)
        self.cropped_labels = self._load_labels(self.labels_dir, self.data_mult)

        self.batches_generate = self._generate_batches(self.cropped_data, self.cropped_labels,
            ~f.border_size: self.border_size, self.inner_size)

        def get_next_batch(self):
            epoch, batchnum, data, labels = self._get_next_batch()

            cropped = self.cropped_data[epoch][batchnum]
            self.__trim_borders(cropped)
            cropped -= self.data_mult

matrix.cpp + (~/Downloads/cuda-convnet/trunk/src/common) - VIM8
void Matrix::updateDims(long int numRows, long int numCols) {
    this->_numRows = numRows;
    this->_numCols = numCols;
    this->_numElements = numRows * numCols;
}

void Matrix::checkBounds(long int startRow, long int endRow, long int startCo...
    ~l, long int endCol) const {
    assert(startRow >= 0 && startRow <= _numRows);
    assert(endRow >= 0 && endRow <= _numRows);
    assert(startCol >= 0 && startCol <= _numCols);
    assert(endCol >= 0 && endCol <= _numCols);
}

/* will return a view if possible */
Matrix& Matrix::slice(long int startRow, long int endRow, long int startCol, l...
    ~ong int endCol) const {
    endRow = endRow < 0 ? this->_numRows : endRow;
    endCol = endCol < 0 ? this->_numCols : endCol;
    _checkBounds(startRow, endRow, startCol, endCol);
    if (!isTrans() && ((startCol == 0 && endCol == this->_numCols) || (startRo...
        ~w == endRow - 1))) {
        return *new Matrix(this->_data + startRow * this->_numCols + startCol,
            ~ endRow - startRow, endCol - startCol);
    } else if (isTrans() && ((startRow == 0 && endRow == this->_numRows) || (s...
        ~tartCol == endCol - 1))) {
        return *new Matrix(this->_data + startCol * this->_numRows + startRow,
            ~ endRow - startRow, endCol - startCol, true);
    }
    Matrix& newSlice = *new Matrix(endRow - startRow, endCol - startCol);
    this->copy(newSlice, startRow, endRow, startCol, endCol, 0, 0);
    return newSlice;
}

convdata.py + (~/Downloads/cuda-convnet/trunk/src/common) - VIM7
class CroppedCIFARDataProvider:
    def __init__(self, data_dir, labels_dir, num_epochs, num_batches_per_epoch, test...
        ~=None, dp_params=None, test_data_dir=None, labels_dir=None, num_epochs=10, num_batches_per_epoch=10,
        ~poch, init_batchnum, dp_params=None):
        self.border_size = dp_params.border_size
        self.inner_size = 32
        self.multiview = dp_params.multiview
        self.num_views = 5*2
        self.data_mult = self.border_size + self.inner_size
        self.num_colors = 3

        for d in self.data_dirs:
            d['data'] = n.recursive_glob(d, ['data'])
            d['labels'] = n.recursive_glob(d, ['labels'])

        self.cropped_data = self._load_data(self.data_dir, self.data_mult)
        self.cropped_labels = self._load_labels(self.labels_dir, self.data_mult)

        self.batches_generate = self._generate_batches(self.cropped_data, self.cropped_labels,
            ~f.border_size: self.border_size, self.inner_size)

        def get_next_batch(self):
            epoch, batchnum, data, labels = self._get_next_batch()

            cropped = self.cropped_data[epoch][batchnum]
            self.__trim_borders(cropped)
            cropped -= self.data_mult

```

```

class AlexNet(nn.Module):
    def __init__(self, num_classes: int = 1000) -> None:
        super(AlexNet, self).__init__()
        self.features = nn.Sequential(
            nn.Conv2d(3, 64, kernel_size=11, stride=4, padding=2),
            nn.ReLU(inplace=True),
            nn.MaxPool2d(kernel_size=3, stride=2),
            nn.Conv2d(64, 192, kernel_size=5, padding=2),
            nn.ReLU(inplace=True),
            nn.MaxPool2d(kernel_size=3, stride=2),
            nn.Conv2d(192, 384, kernel_size=3, padding=1),
            nn.ReLU(inplace=True),
            nn.Conv2d(384, 256, kernel_size=3, padding=1),
            nn.ReLU(inplace=True),
            nn.Conv2d(256, 256, kernel_size=3, padding=1),
            nn.ReLU(inplace=True),
            nn.MaxPool2d(kernel_size=3, stride=2),
        )
        self.avgpool = nn.AdaptiveAvgPool2d((6, 6))
        self.classifier = nn.Sequential(
            nn.Dropout(),
            nn.Linear(256 * 6 * 6, 4096),
            nn.ReLU(inplace=True),
            nn.Dropout(),
            nn.Linear(4096, 4096),
            nn.ReLU(inplace=True),
            nn.Linear(4096, num_classes),
        )

    def forward(self, x: torch.Tensor) -> torch.Tensor:
        x = self.features(x)
        x = self.avgpool(x)
        x = torch.flatten(x, 1)
        x = self.classifier(x)
        return x

```

AlexNet (2012) CUDA: 12k lines Python: 2k lines C++: 1k lines

AlexNet (2024) PyTorch: 35 lines



assembly :



:: CUDA :



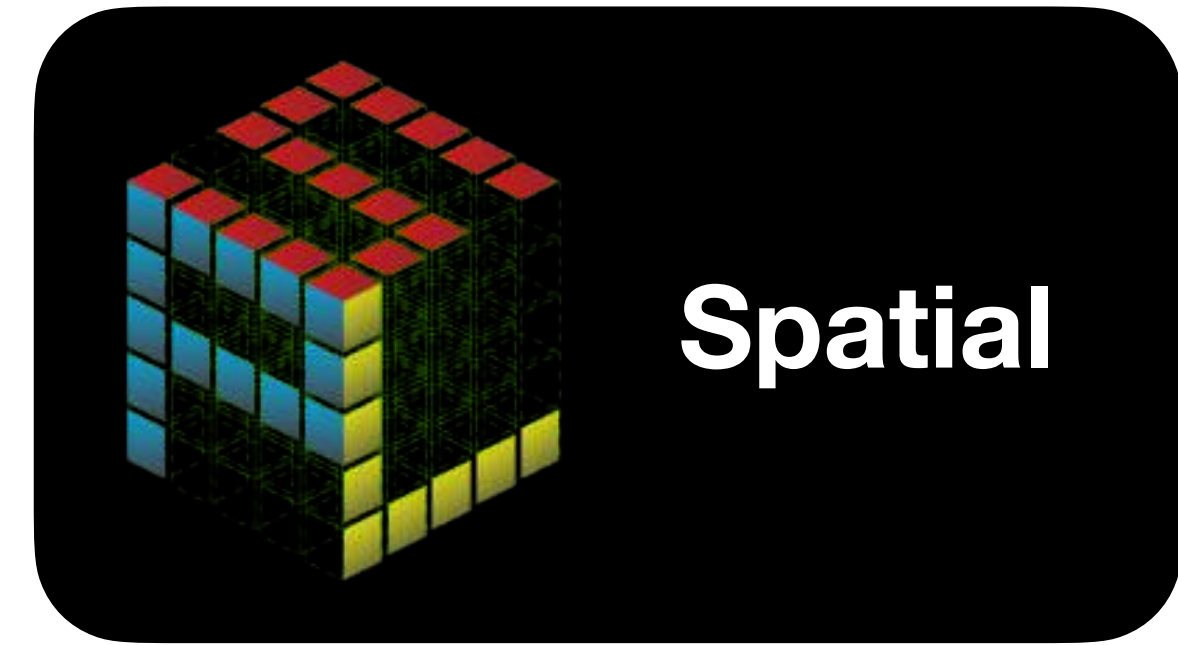
assembly :



::

HDLs :

**Accelerator
Design
Languages**



[Koeplinger et al., PLDI 2018]

assembly :



::

HDLs :



eteroCL

[Lai et al., FPGA 2019]

Aetherling

[Durst et al., PLDI 2020]

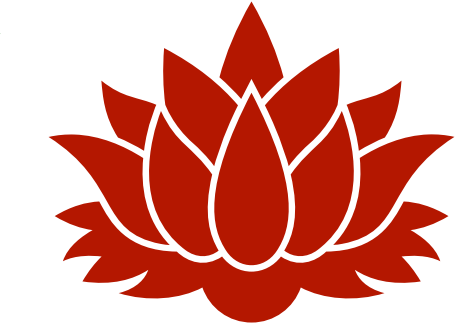
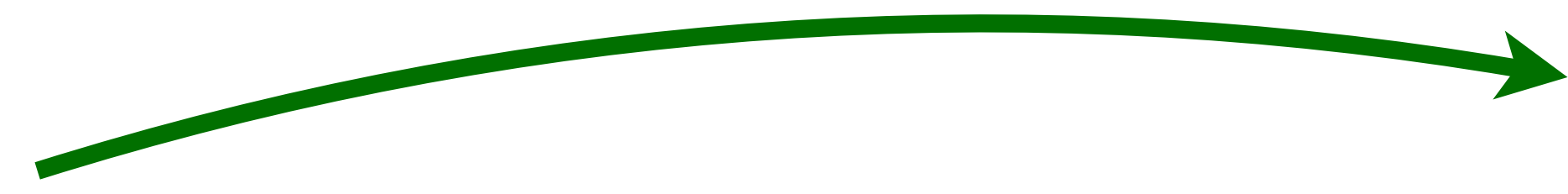
Dahlia



[Nigam et al., PLDI 2020]

Academic ADLs

debuggers [Berlstein et al., ASPLOS 2023]
profilers
module systems



assembly :



:: HDLs :

**rethinking the system stack
for reconfigurable hardware (FPGAs)**
[Vega et al., PLDI 2021]

**type systems for
modular reasoning**
[Nigam et al., PLDI 2023]

1. NEVER LISTEN TO ADVICE
2. ALWAYS LISTEN TO ADVICE
3. IT'S NOT NAÏVETÉ; IT'S BEGINNER'S MIND
4. THE BEST COMPUTER ARCHITECTURE RESEARCH ALWAYS "CHEATS"
5. DON'T DO FRAMEWORKS
6. START A BLOG
7. TALK TO YOUR FRIENDS