

GPGPU Performance and Power Estimation Using Machine Learning

Gene Wu – UT Austin

Joseph Greathouse – AMD Research

Alexander Lyashevsky – AMD Research

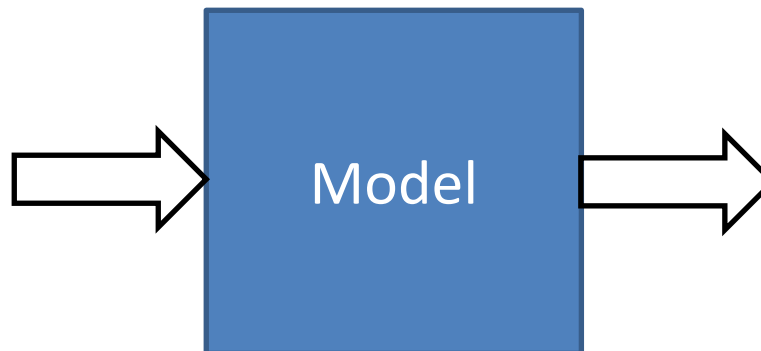
Nuwan Jayasena – AMD Research

Derek Chiou – UT Austin

Goals

- Create GPU power & performance scalability models:
 - Capable of predicting for a wide range of settings
 - Number of Compute Units(CUs) or parallel cores
 - GPU Core(Engine) Frequency
 - Memory Frequency
 - Predict many hardware configurations from data gathered on a single configuration

- **Compute Units = 4**
- Execution Time/Power
- Performance Counter Values



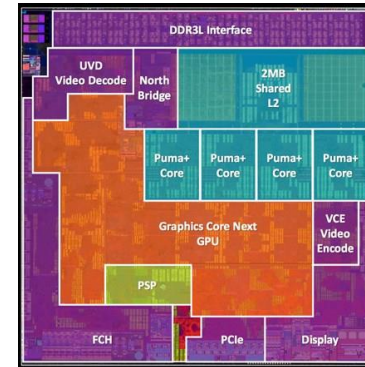
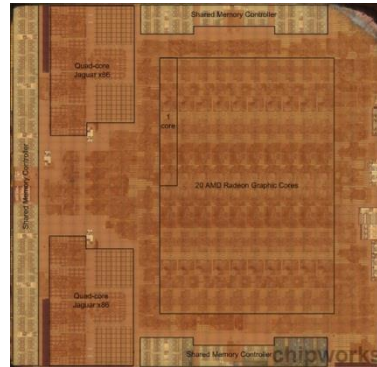
- **Compute Units = 32**
- Predicted Execution Time/Power

Why Power and Performance Estimation?

- Feedback to programmer.

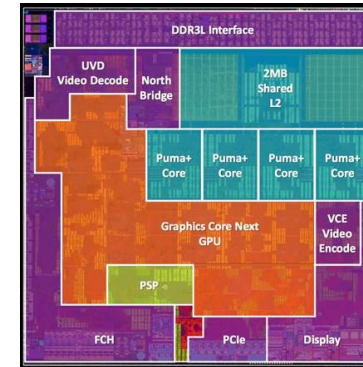
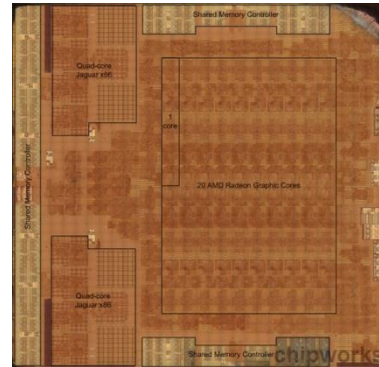
Why Power and Performance Estimation?

- Feedback to programmer.
- HW Design Exploration (e.g. Semi-Custom)

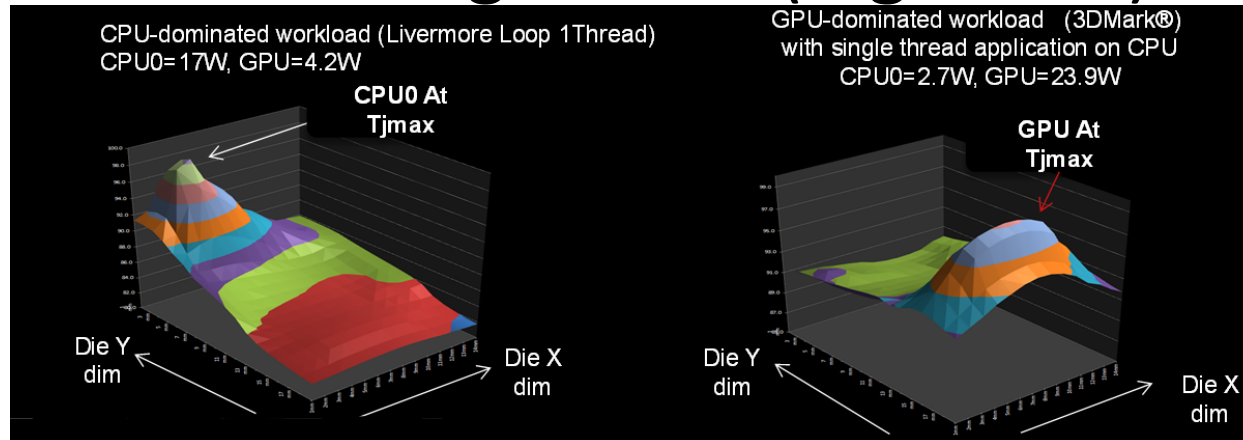


Why Power and Performance Estimation?

- Feedback to programmer.
- HW Design Exploration (e.g. Semi-Custom)



- Online reconfiguration (e.g. DVFS)

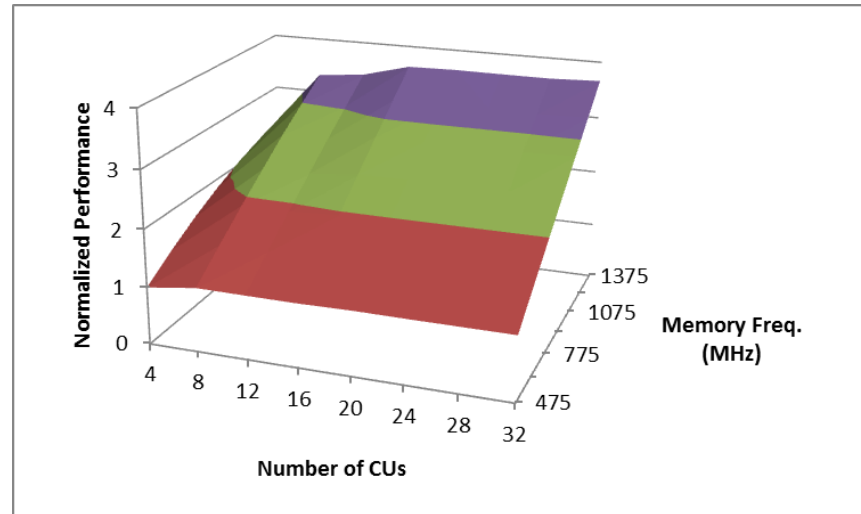


Outline

- Goals
- Model Overview
- Model Construction
- Results

Base to Target Config. Execution

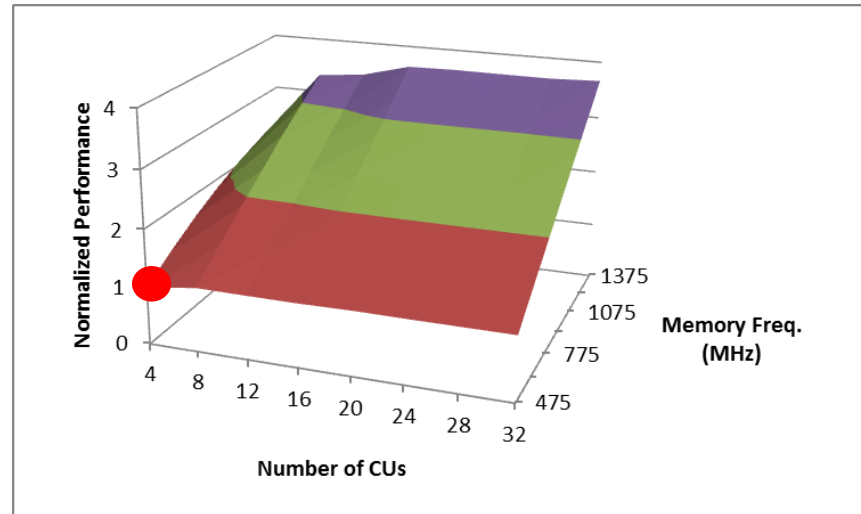
- Hardware Configuration
 - Compute unit(CU) count
 - Engine frequency
 - Memory frequency



- The hardware configuration from which measurements are taken is the **Base Hardware Configuration**
- The hardware configuration that we wish to predict performance/power at is the **Target Hardware Configuration**

Base to Target Config. Execution

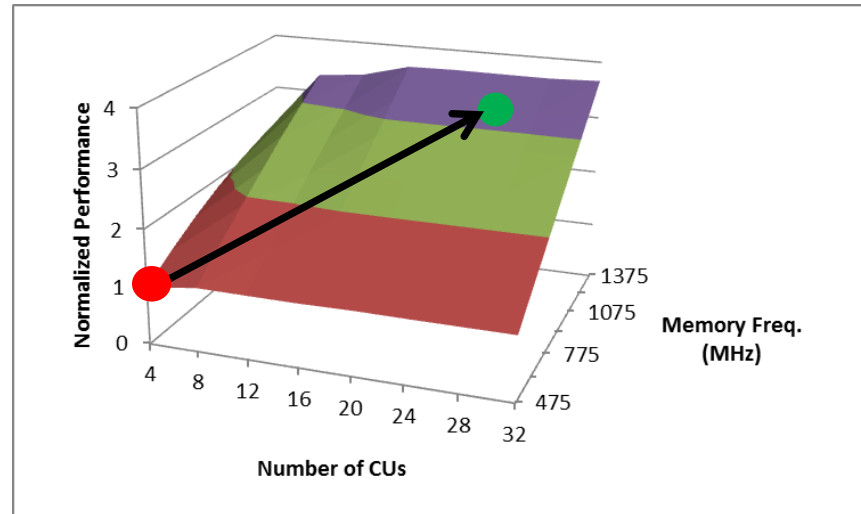
- Hardware Configuration
 - Compute unit(CU) count
 - Engine frequency
 - Memory frequency



- The hardware configuration from which measurements are taken is the **Base Hardware Configuration**
- The hardware configuration that we wish to predict performance/power at is the **Target Hardware Configuration**

Base to Target Config. Execution

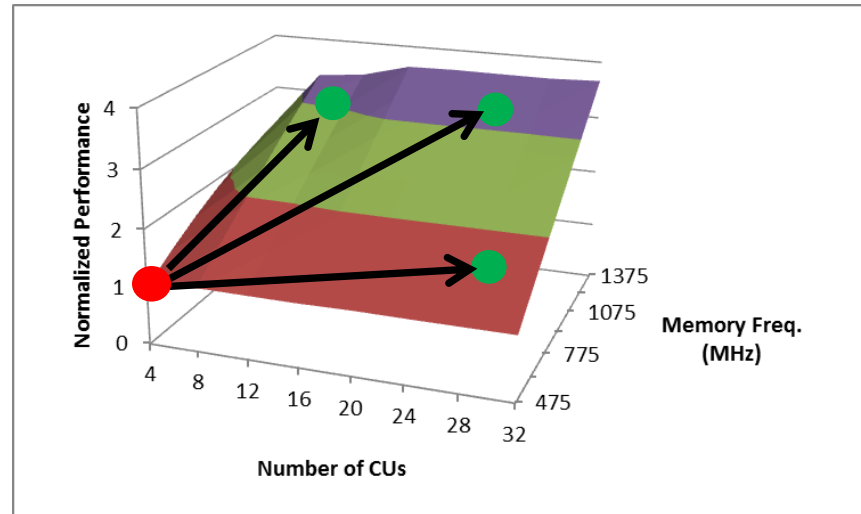
- Hardware Configuration
 - Compute unit(CU) count
 - Engine frequency
 - Memory frequency



- The hardware configuration from which measurements are taken is the **Base Hardware Configuration**
- The hardware configuration that we wish to predict performance/power at is the **Target Hardware Configuration**

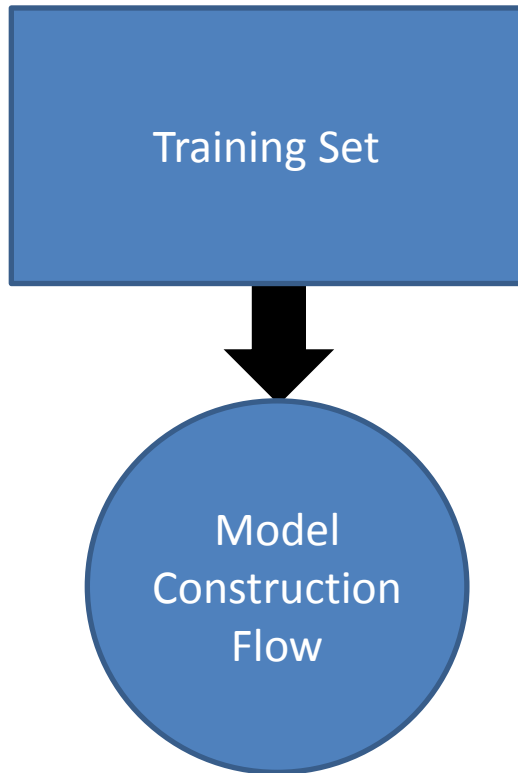
Base to Target Config. Execution

- Hardware Configuration
 - Compute unit(CU) count
 - Engine frequency
 - Memory frequency

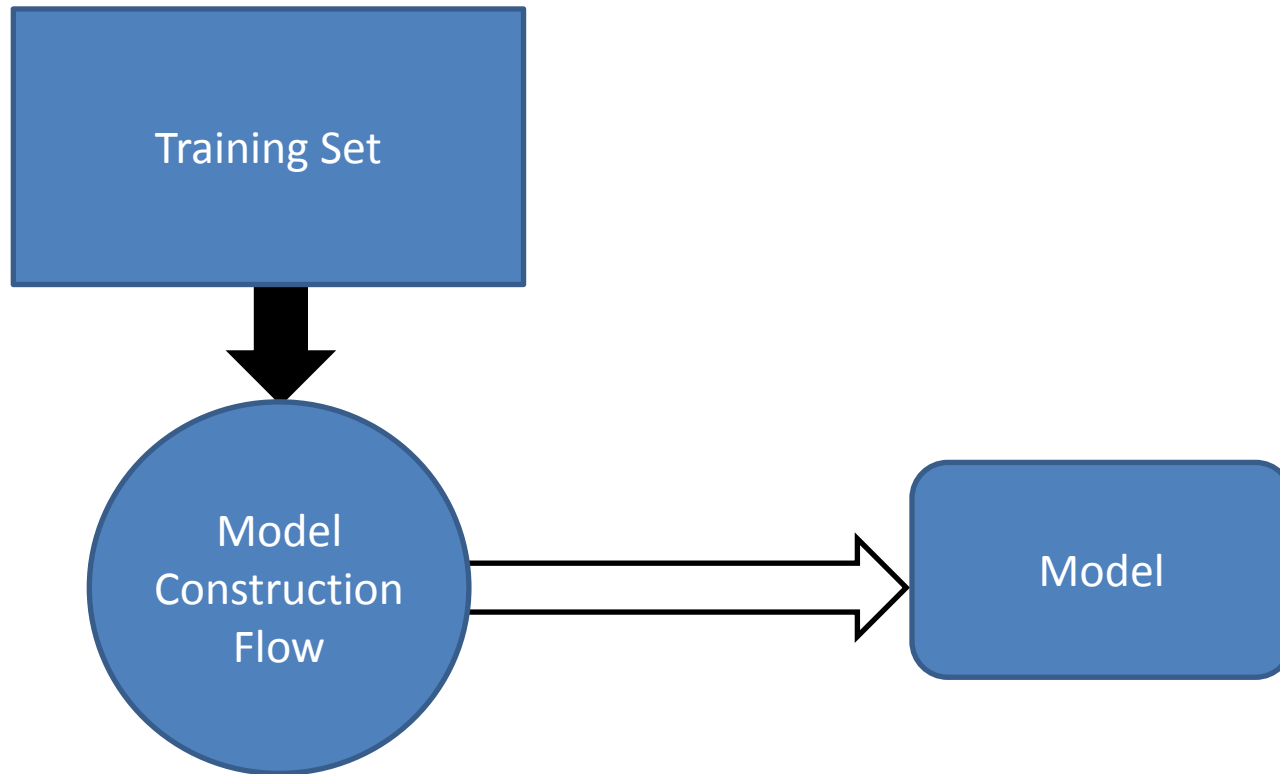


- The hardware configuration from which measurements are taken is the **Base Hardware Configuration**
- The hardware configuration that we wish to predict performance/power at is the **Target Hardware Configuration**

Model Construction and Usage Flow

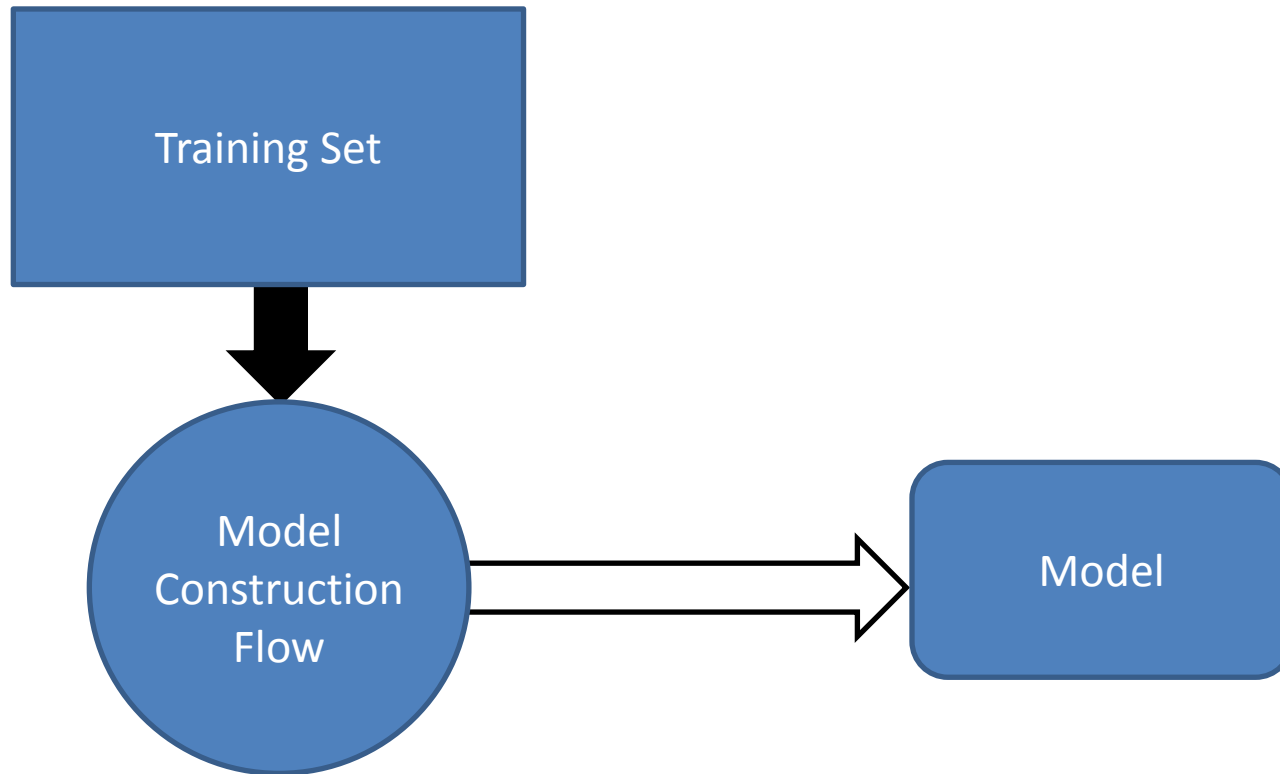


Model Construction and Usage Flow

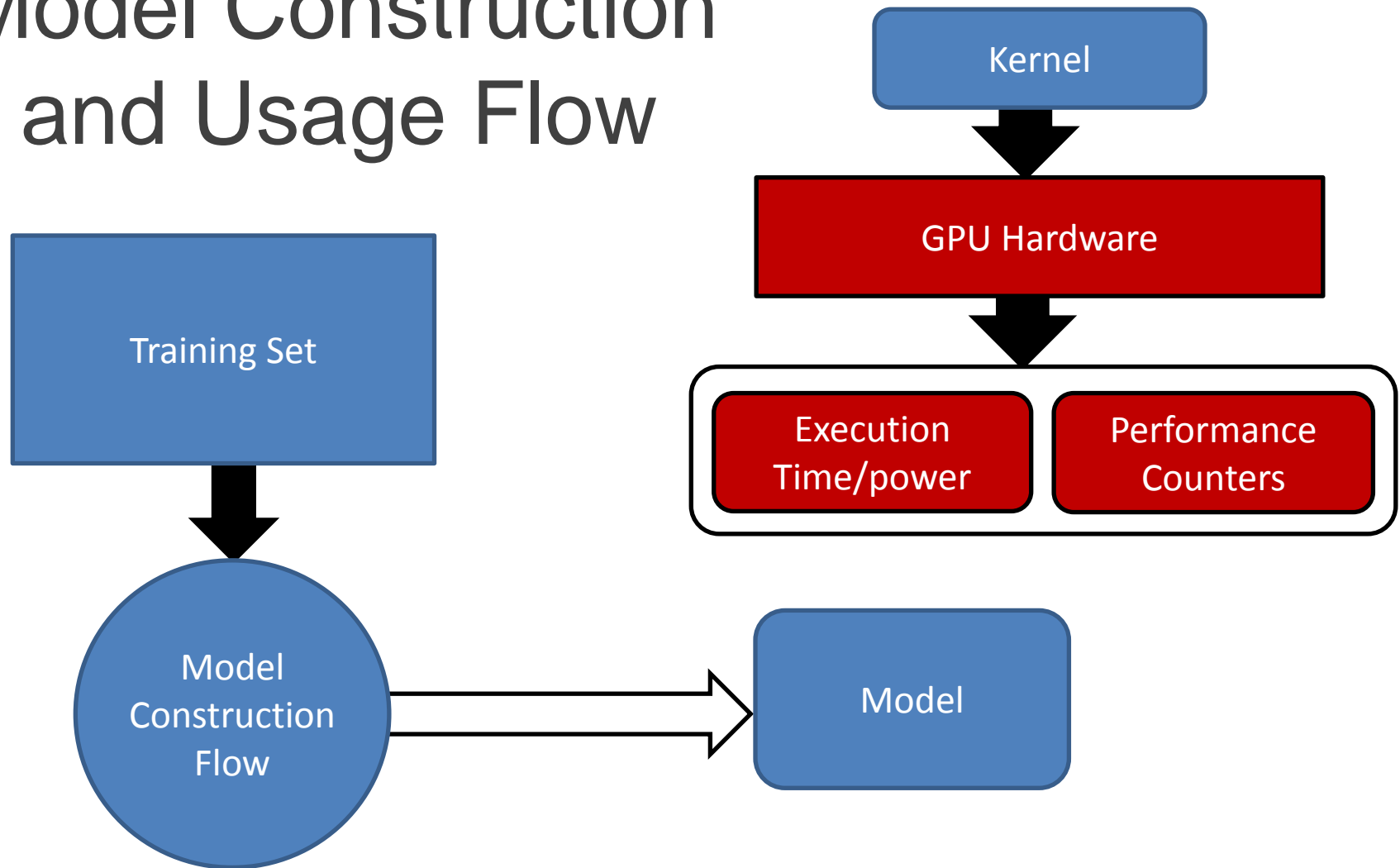


Model Construction and Usage Flow

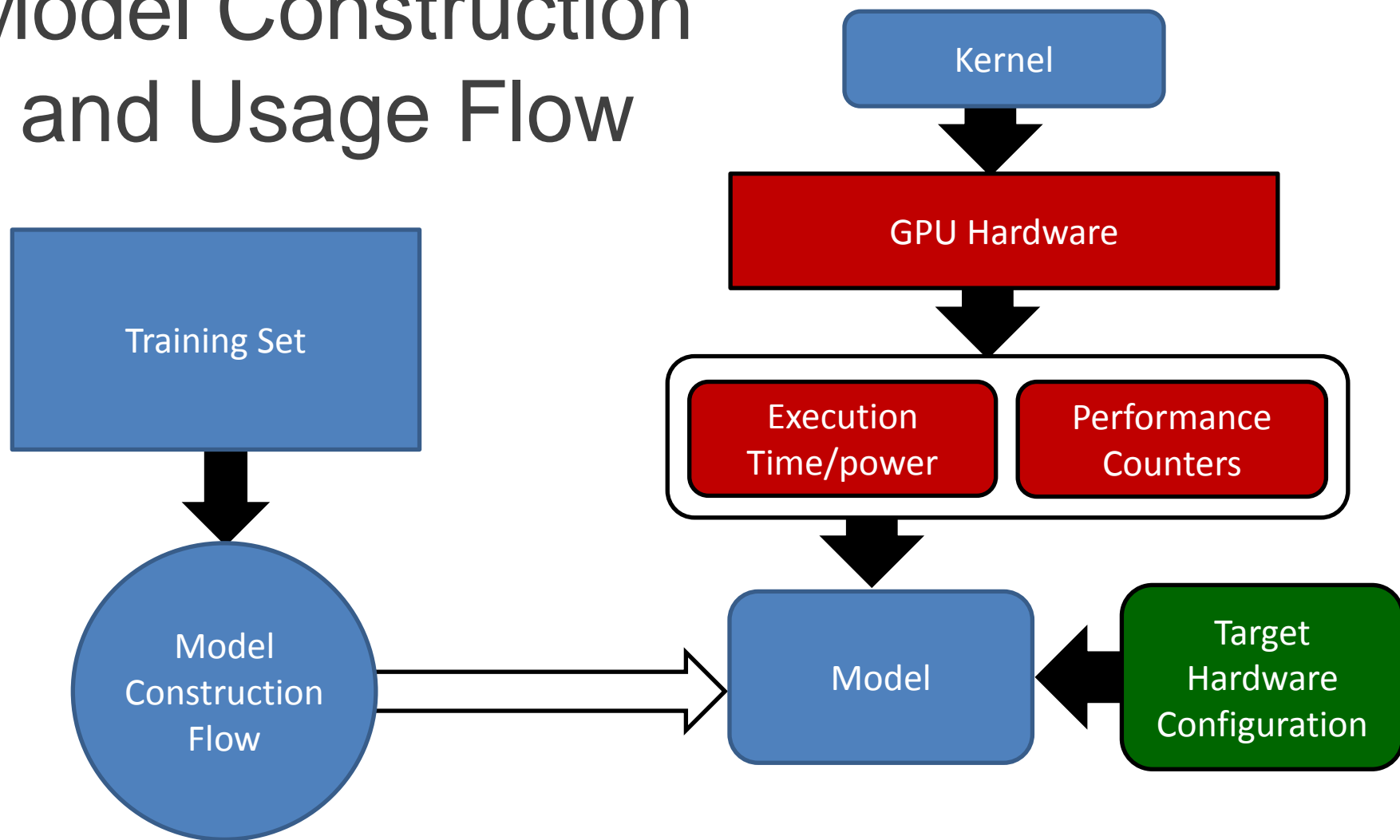
Kernel



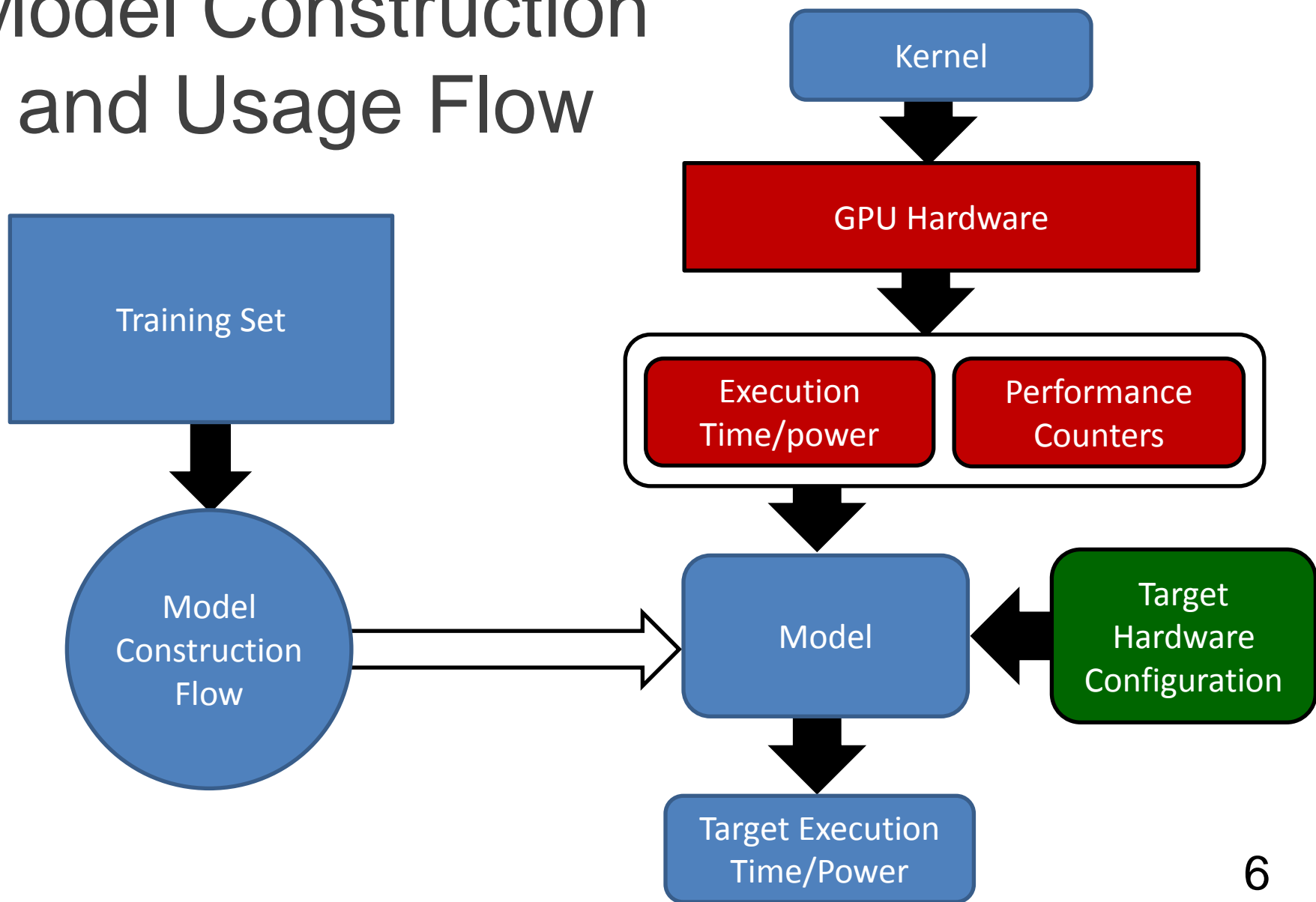
Model Construction and Usage Flow



Model Construction and Usage Flow



Model Construction and Usage Flow



Training Set

CU count, Engine freq., Mem. Freq.



Kernel name	4,300,375	8,300,375	...	32,1000,1375	Perf. Count 1.	Perf. Count. 2	...
Kernel 1							
Kernel 2							
.....							
Kernel N							

Execution Times/Power

Performance Counter Values
gathered on base hardware
configuration

Outline

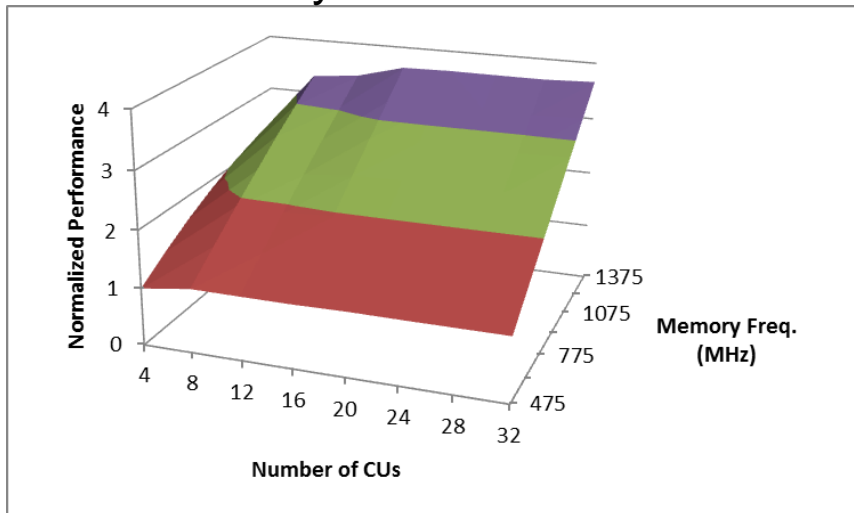
- Goals
- Model Overview
- Model Construction
- Results

Model Construction

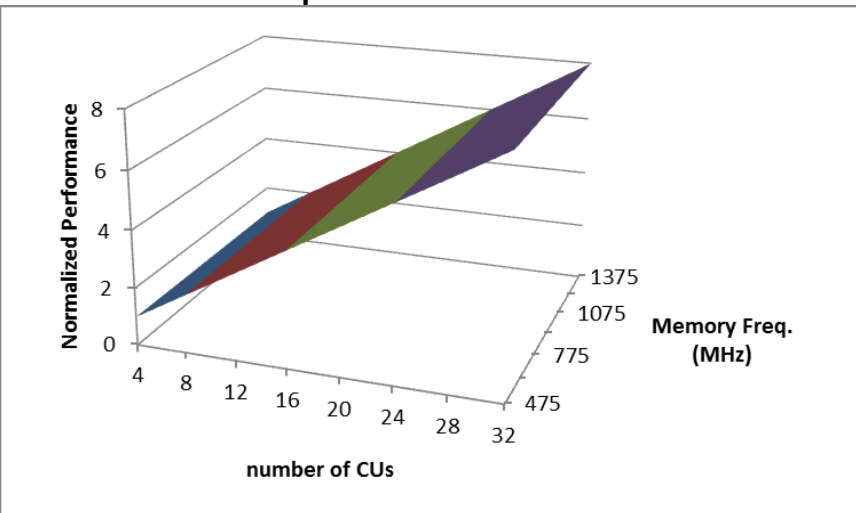
- Phase 1: Form clusters of training kernels that scale similarly
- Phase 2: Build a classifier to map kernel performance counter values to specific clusters

Kernel Scaling Behaviors

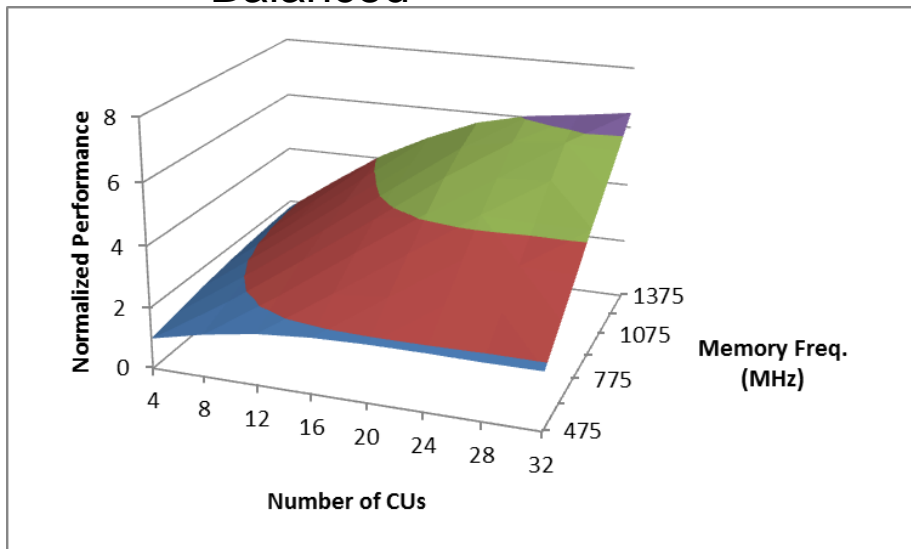
Memory Bound



Compute Bound



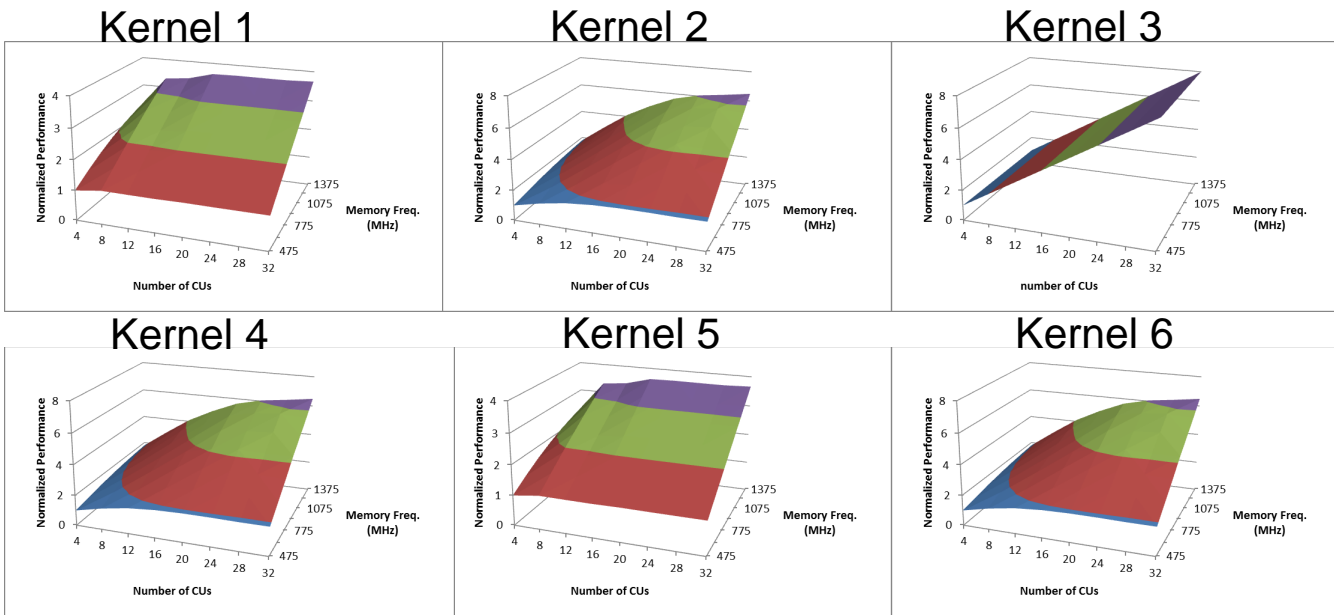
Balanced



- Found many other patterns during this study

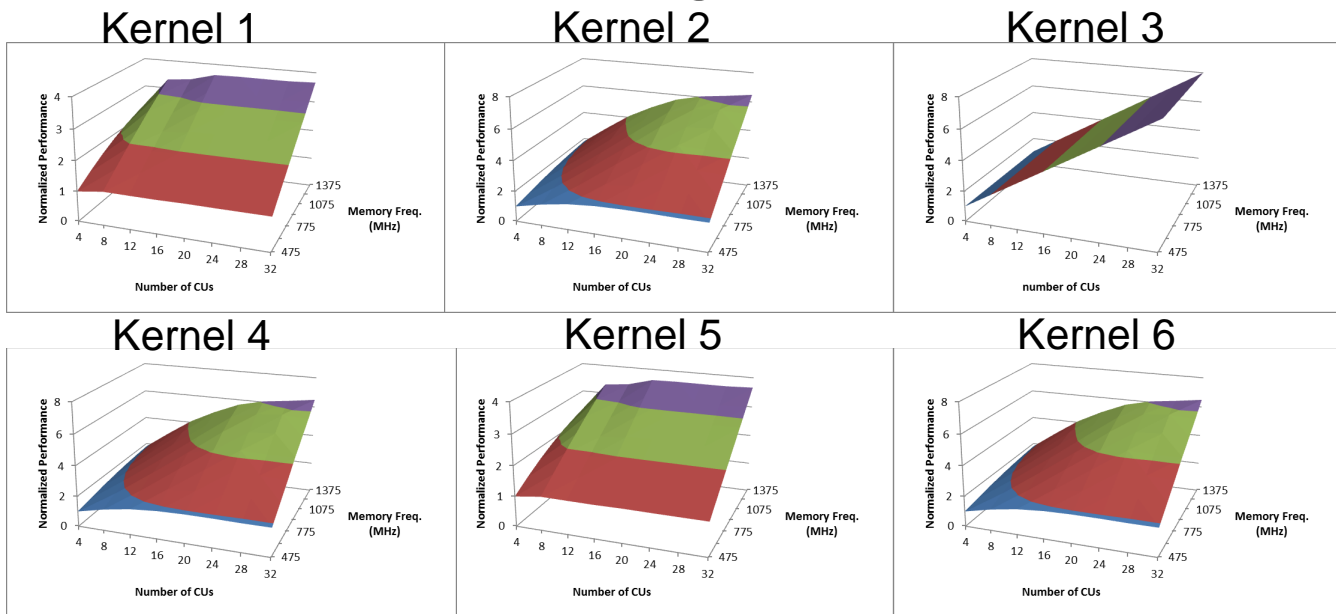
Phase 1: Clustering

Training Set



Phase 1: Clustering

Training Set



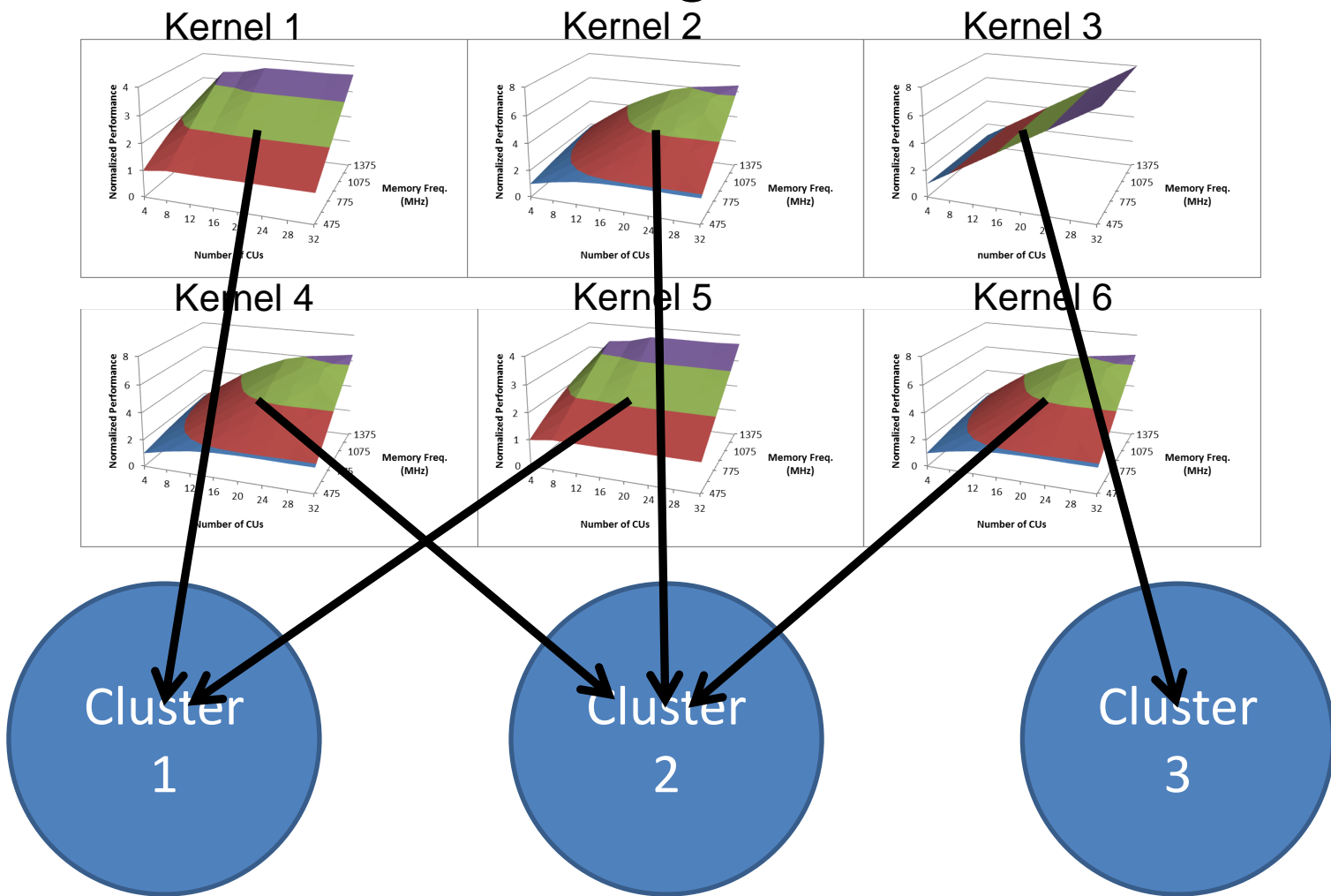
Cluster
1

Cluster
2

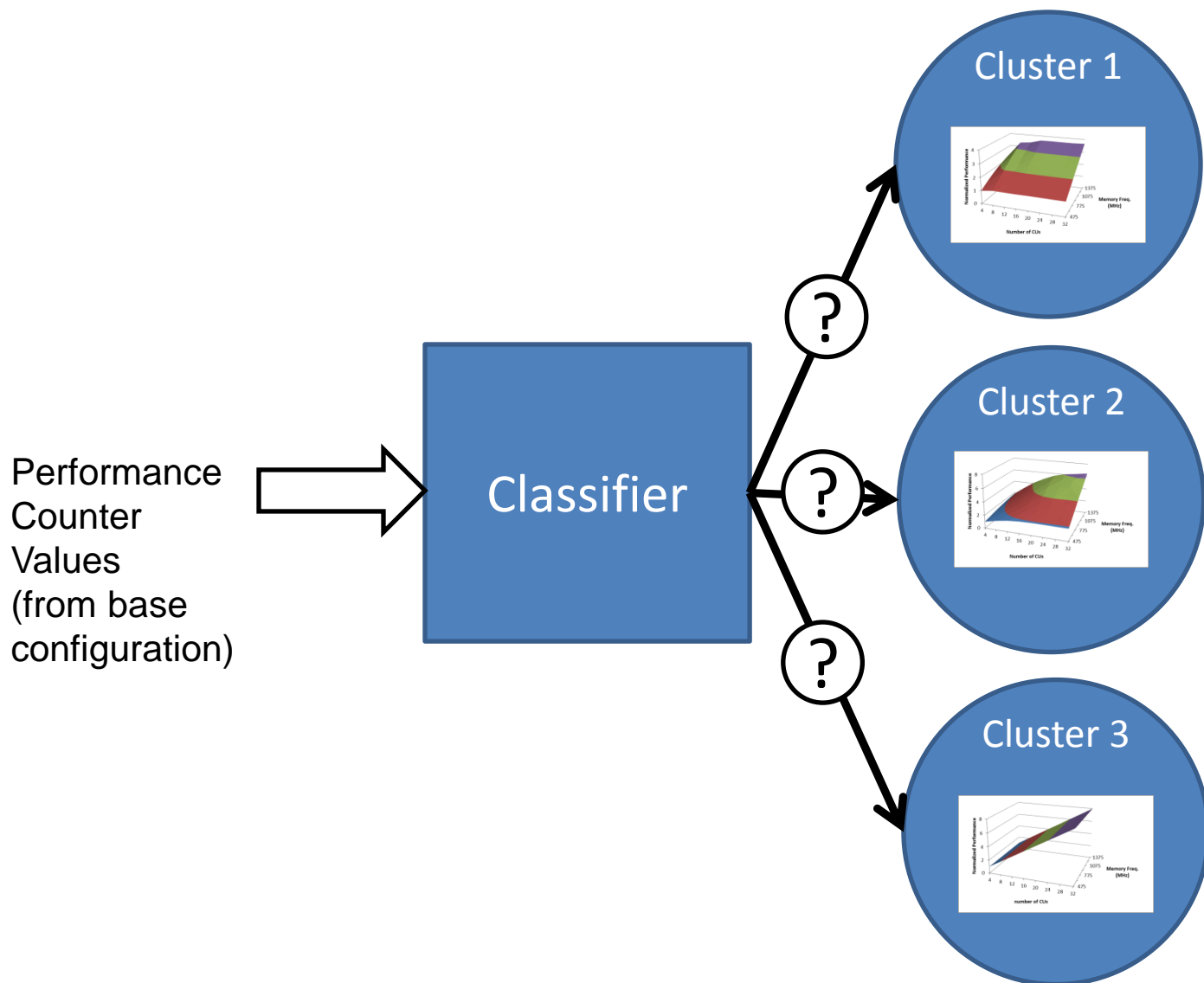
Cluster
3

Phase 1: Clustering

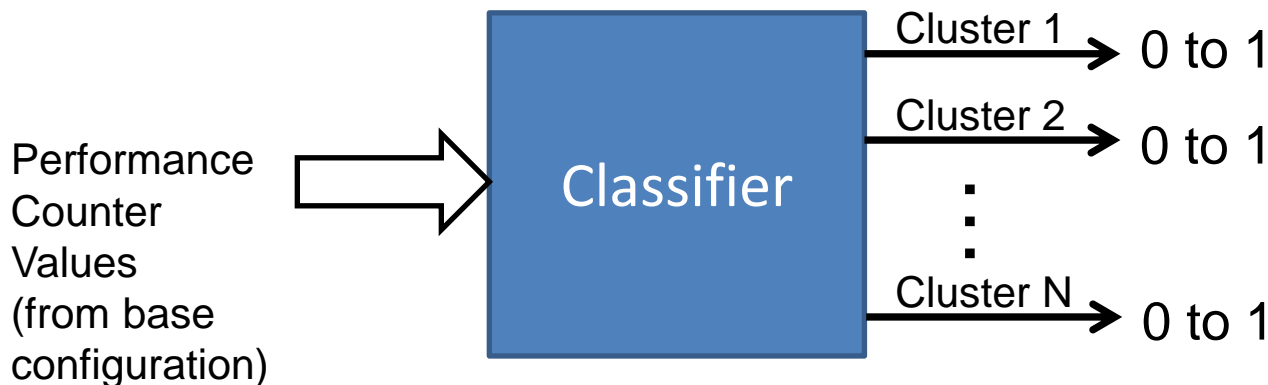
Training Set



Phase 2: Classification



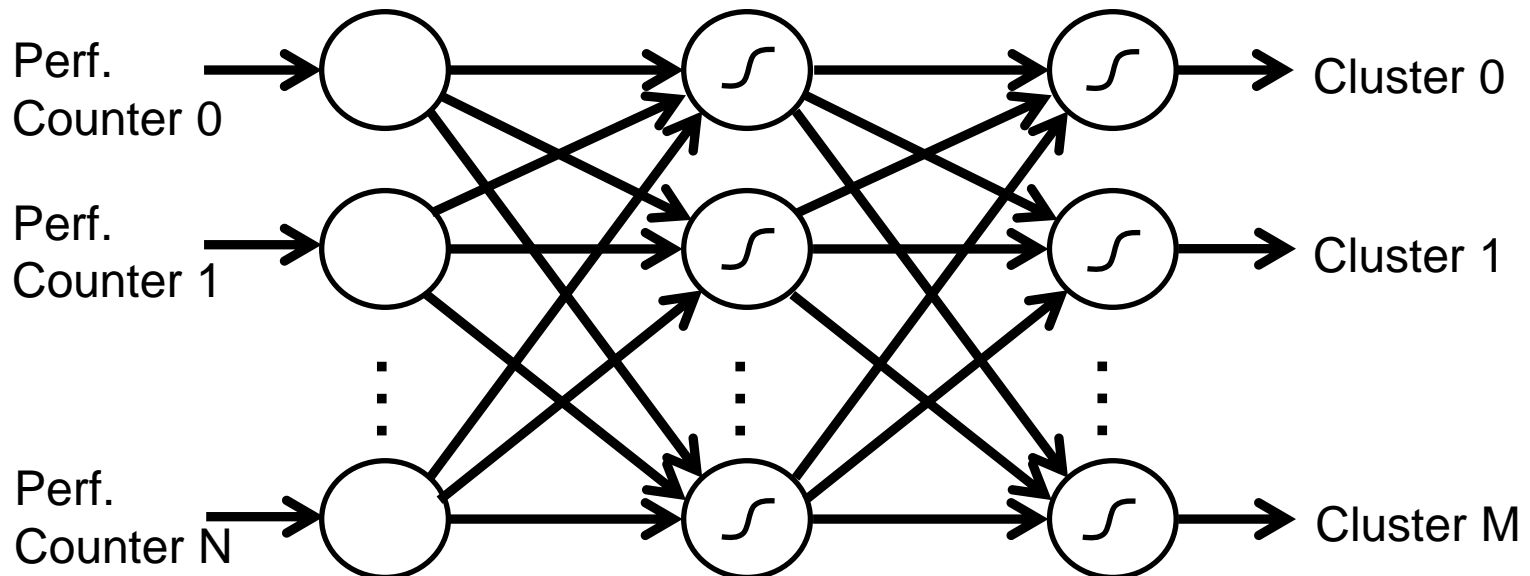
Classifier



- Inputs:
 - Performance counter values
- Outputs:
 - One output per cluster
 - Output values between 0 and 1
 - Cluster with highest output is chosen
 - Ideally a one hot encoding at outputs

Classifier: Neural Network Topology

- 3 layer, fully connected network
 - Input layer: linear
 - Number of neurons equals number of features
 - Hidden layer: sigmoid
 - Number of neurons equals number of clusters
 - Output layer: sigmoid
 - Number of neurons equals number of clusters



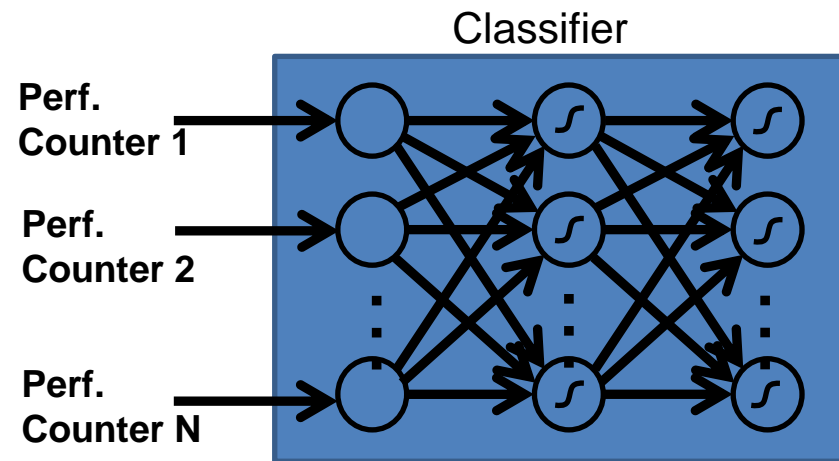
Putting It All Together

**Perf.
Counter 1**

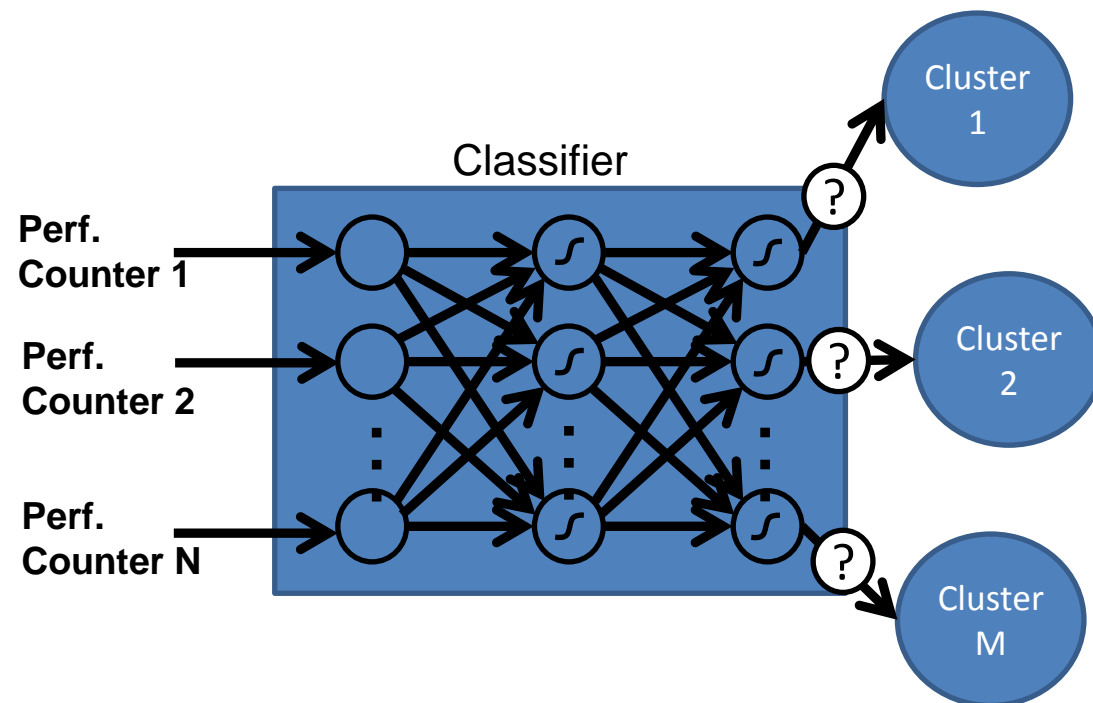
**Perf.
Counter 2**

**Perf.
Counter N**

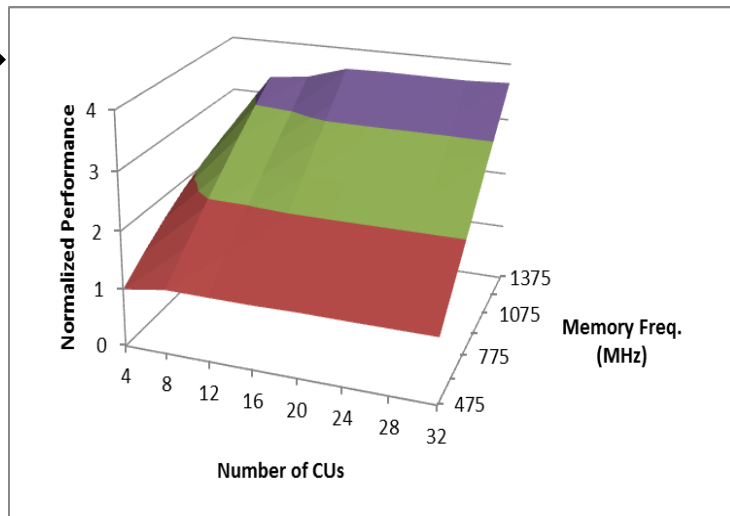
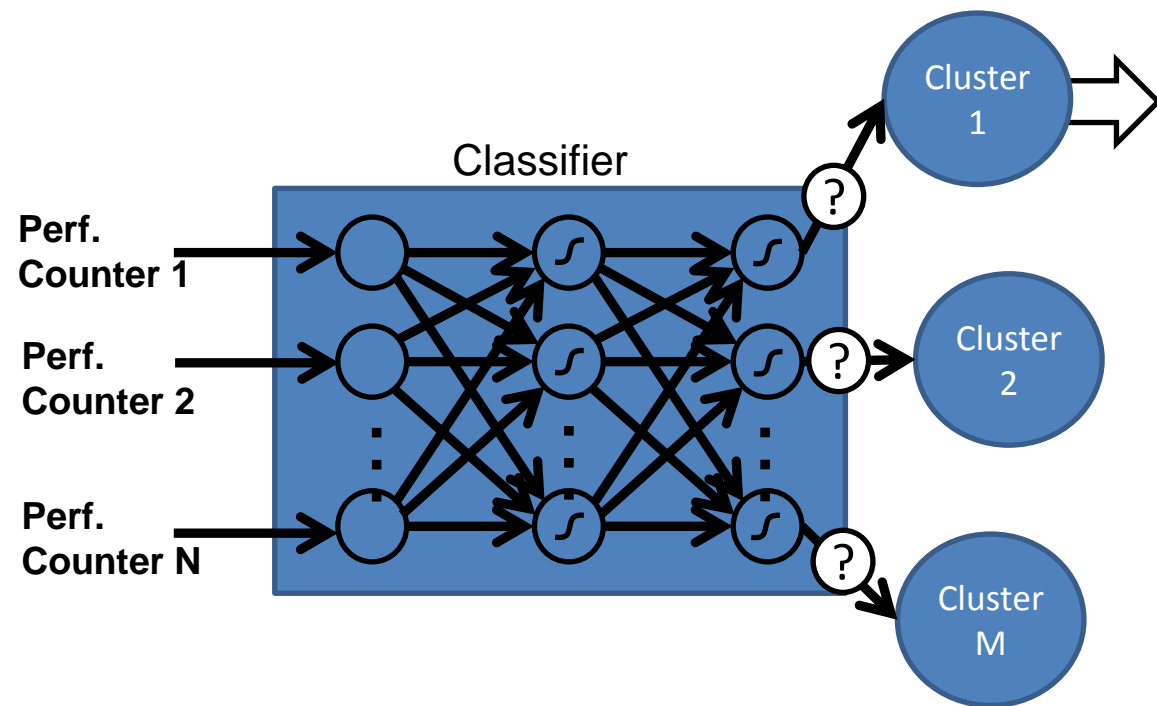
Putting It All Together



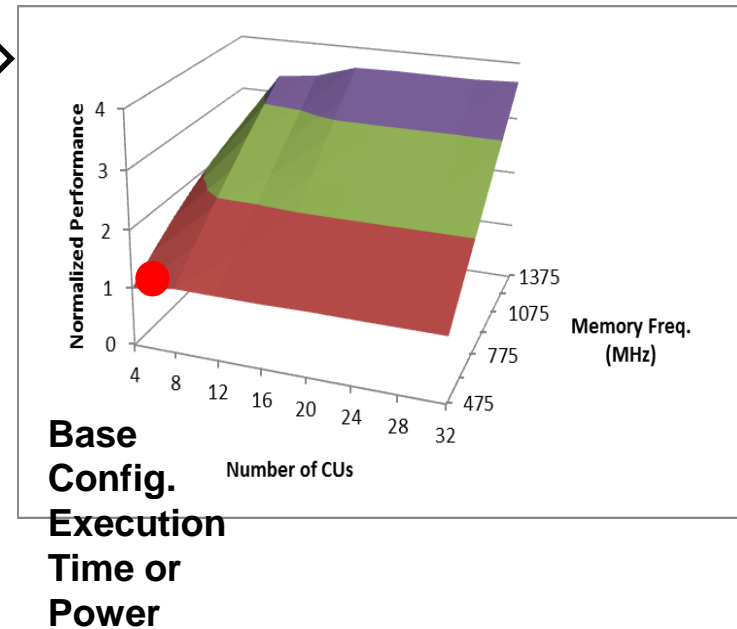
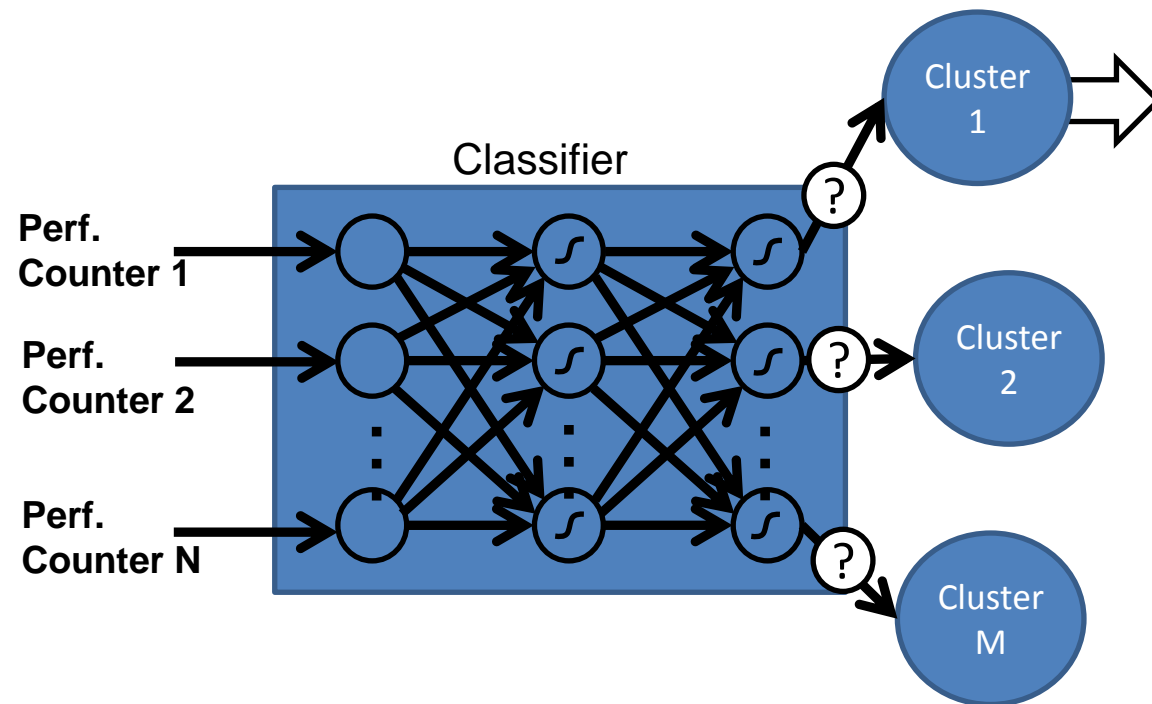
Putting It All Together



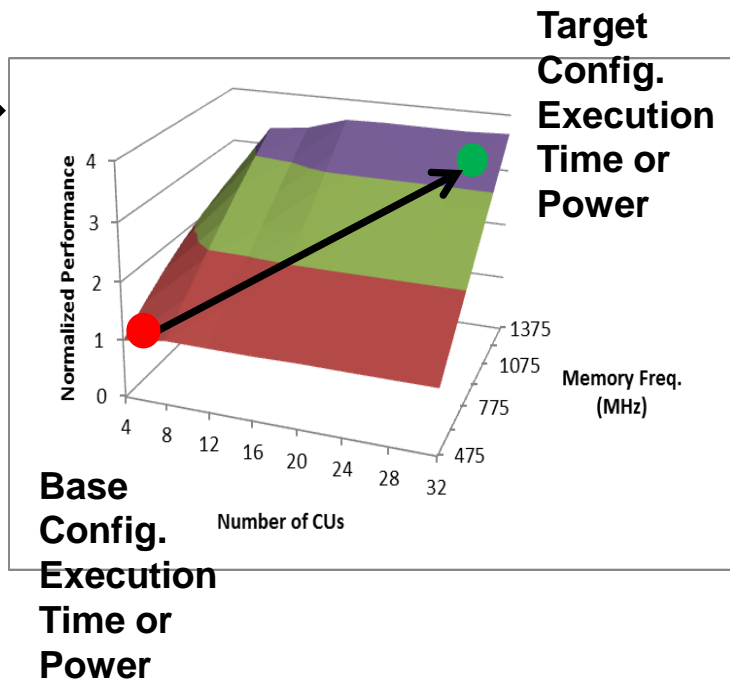
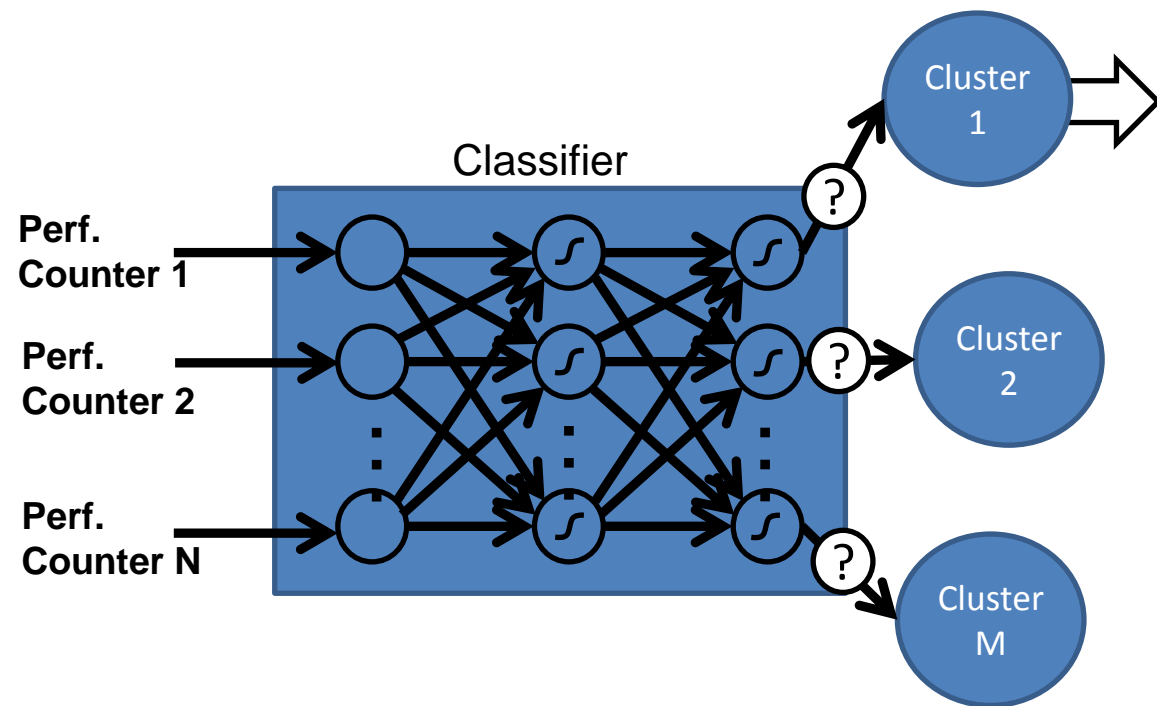
Putting It All Together



Putting It All Together



Putting It All Together



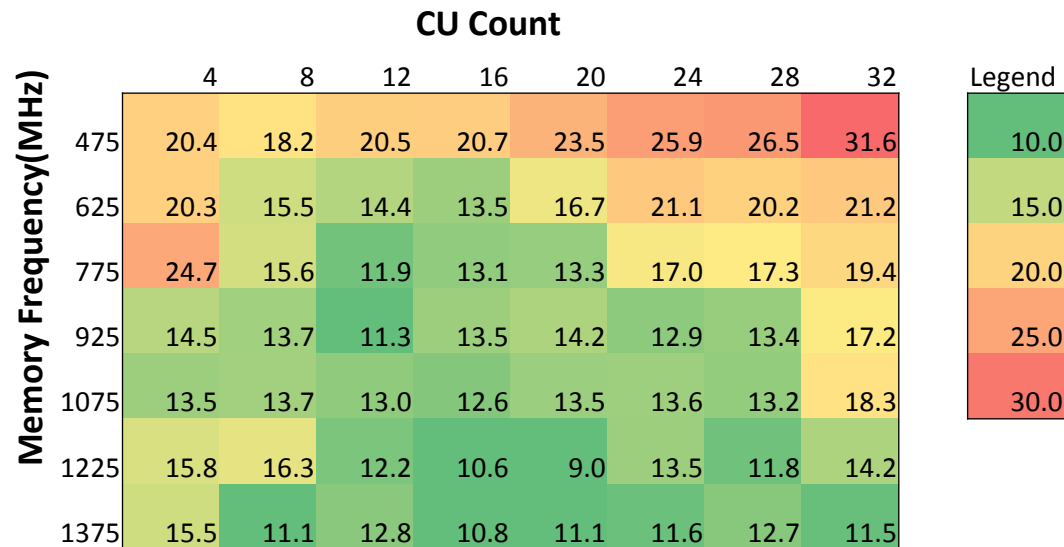
Outline

- Goals
- Model Overview
- Model Construction
- Results

Experimental Setup

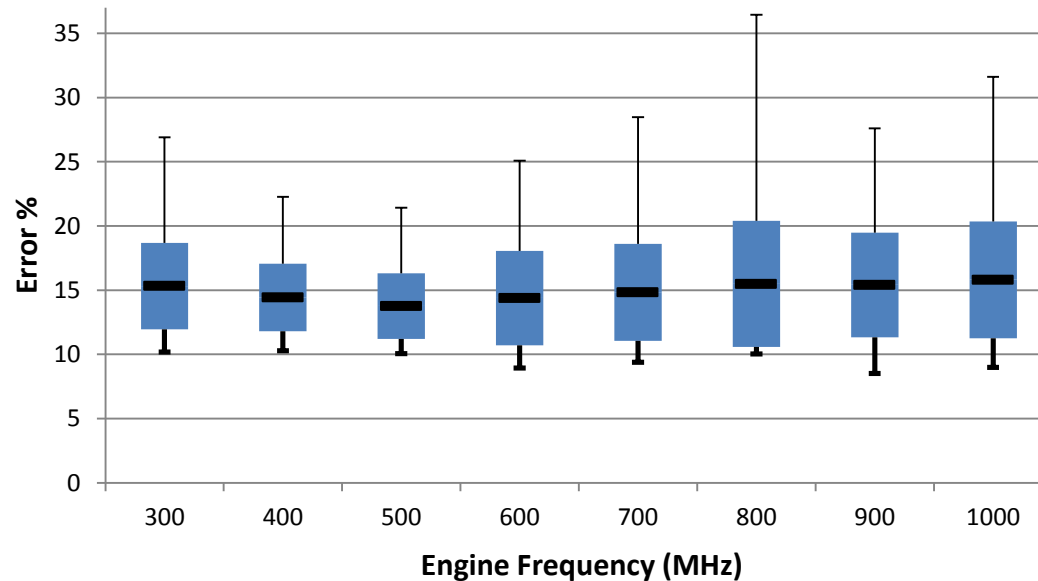
- Measurements gathered on a AMD Radeon HD 7970 GPU
- 8 CU settings:
 - 4, 8, 12, 16, 20, 24, 28, 32
- 8 Engine Frequencies:
 - 300, 400, 500, 600, 700, 800, 900, 1000 (MHz)
- 7 Memory Frequencies:
 - 475, 625, 775, 925, 1075, 1225, 1375 (MHz)
- 448(8x8x7) possible hardware configurations
- 108 OpenCL kernels:
 - 86 kernels (80%) for training
 - 22 kernels (20%) for validation

Accuracy vs. Base Configuration



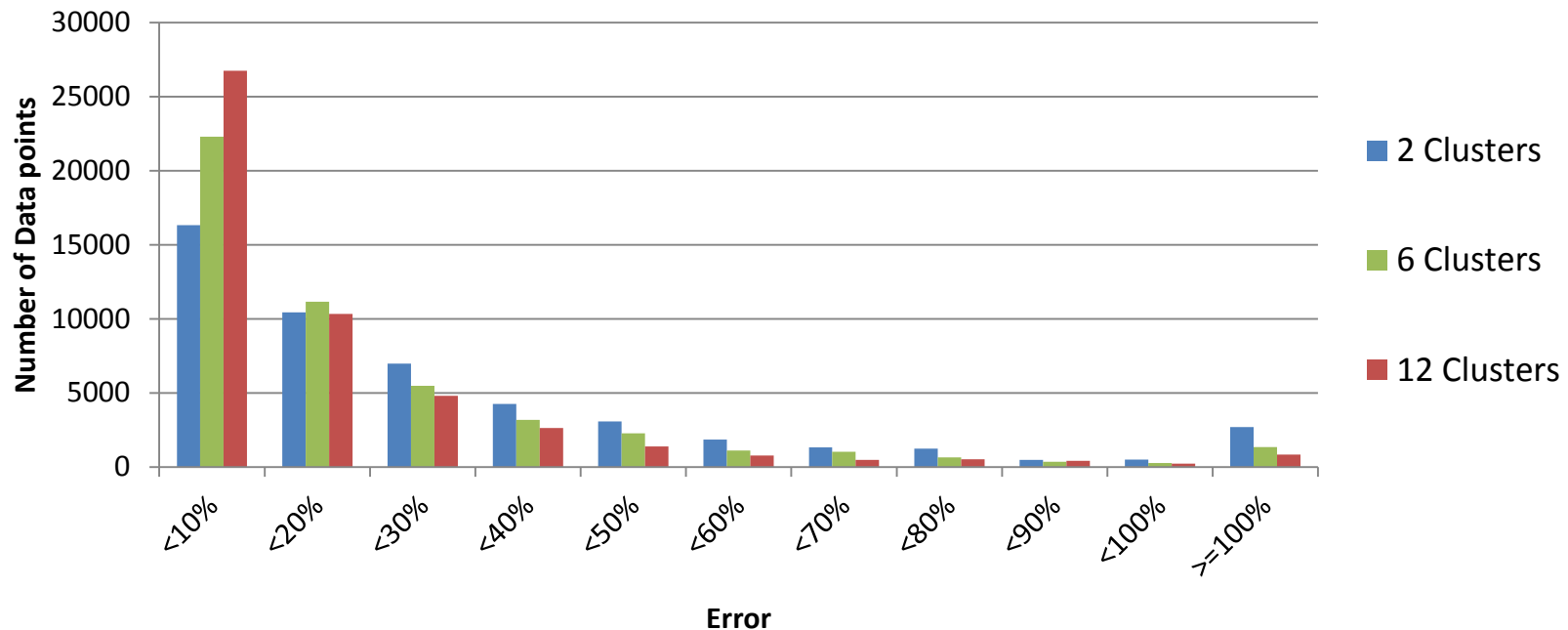
- Base configuration engine frequency fixed at 1000 MHz
- 12 Clusters
- Each entry is the average error of all validation kernels on all 447 possible target configurations (22 kernels x 447 target configs = 9834 predictions)
- Error higher when base configurations has an unbalanced compute to bandwidth ratio

Accuracy vs. Base Configuration



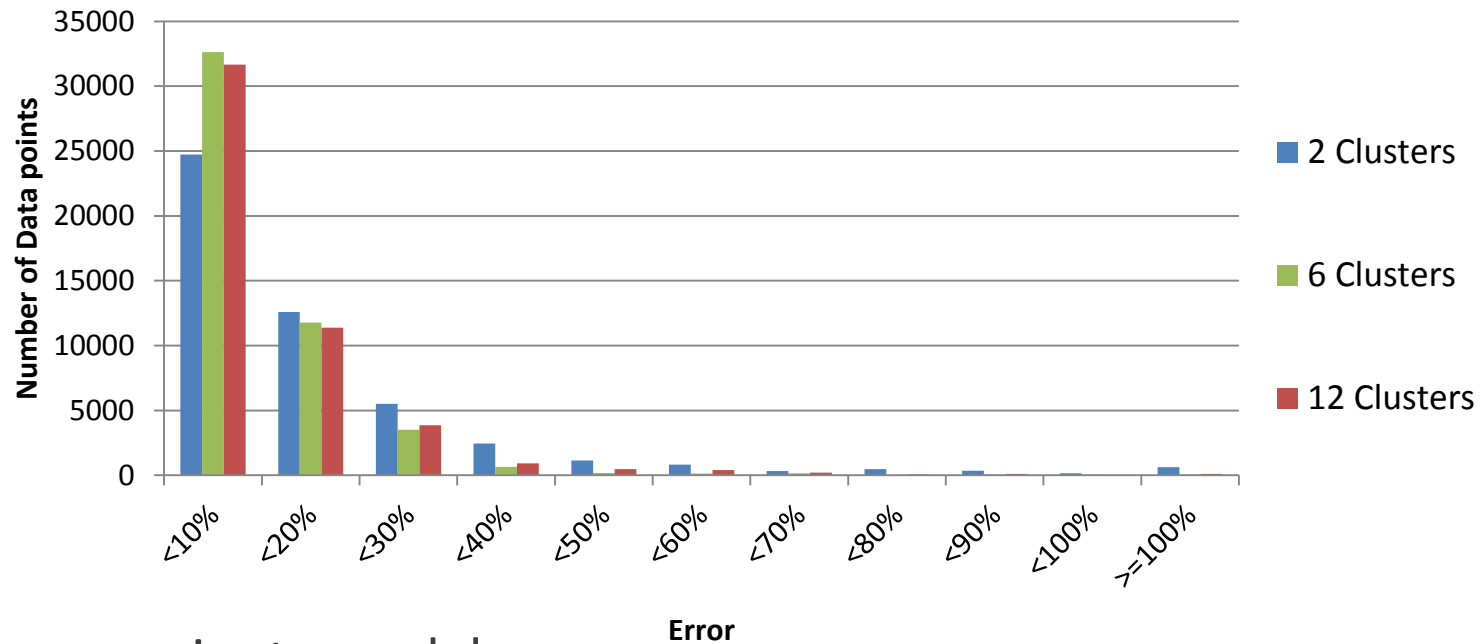
- Lowest error at 500 MHz engine frequency
 - Avg: 13.7%
 - Standard deviation: 2.6%
 - Max: 21.4%
 - Min: 10.1%

Performance Error Distribution



- 447 target configurations
- 22 validation kernels
- 5 base configurations:
 - 32.300.475, 32.300.1375, 32.700.925, 32.1000.475, 32.1000.1375
- $447 \times 22 \times 5 = 49170$ total data points

Power Error Distribution



- Power easier to model
- Modeling a model
- Average Error:
 - 2 Clusters: 11.4%
 - 6 Clusters: 9.1%
 - 12 Clusters: 10.1%

Summary

- GPU power and performance models
 - Constructed with K-means clustering and neural networks
- Performance model average error:
 - Around 10% for the best base hardware configurations
- Power model average error:
 - Around 10%
- Less than a millisecond for each prediction

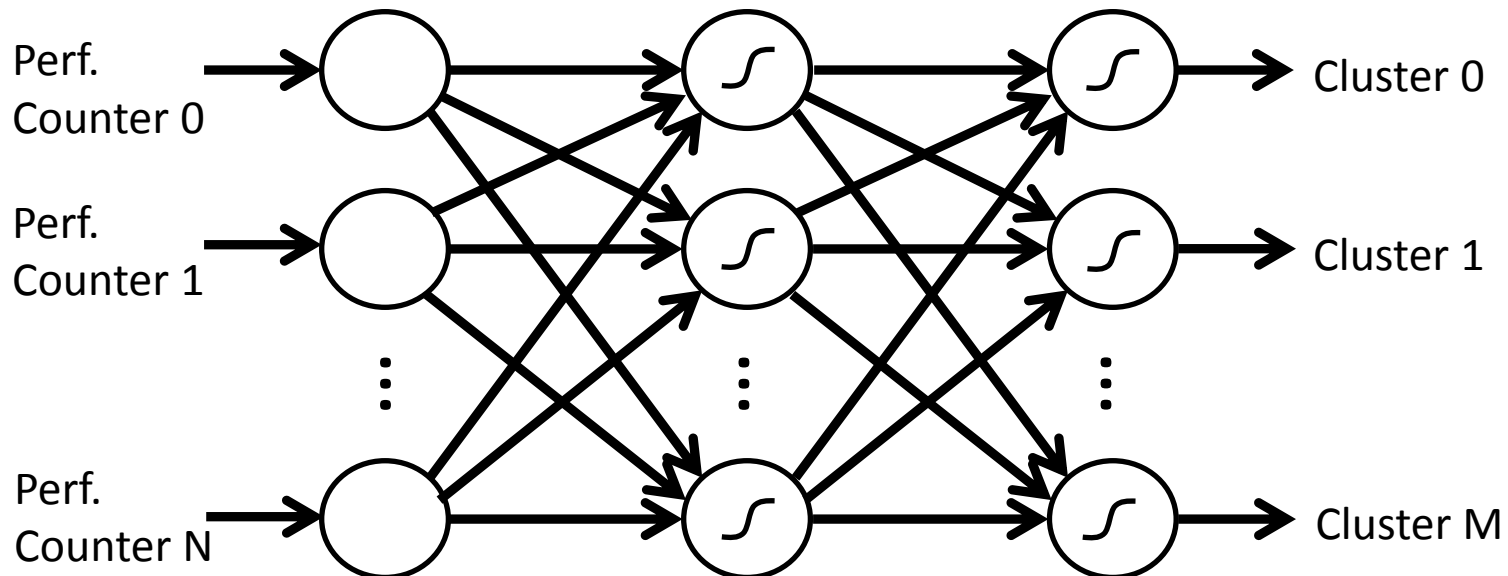
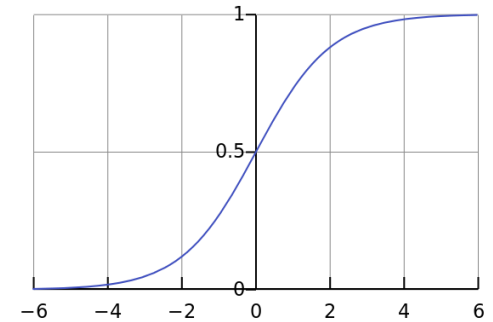
Questions



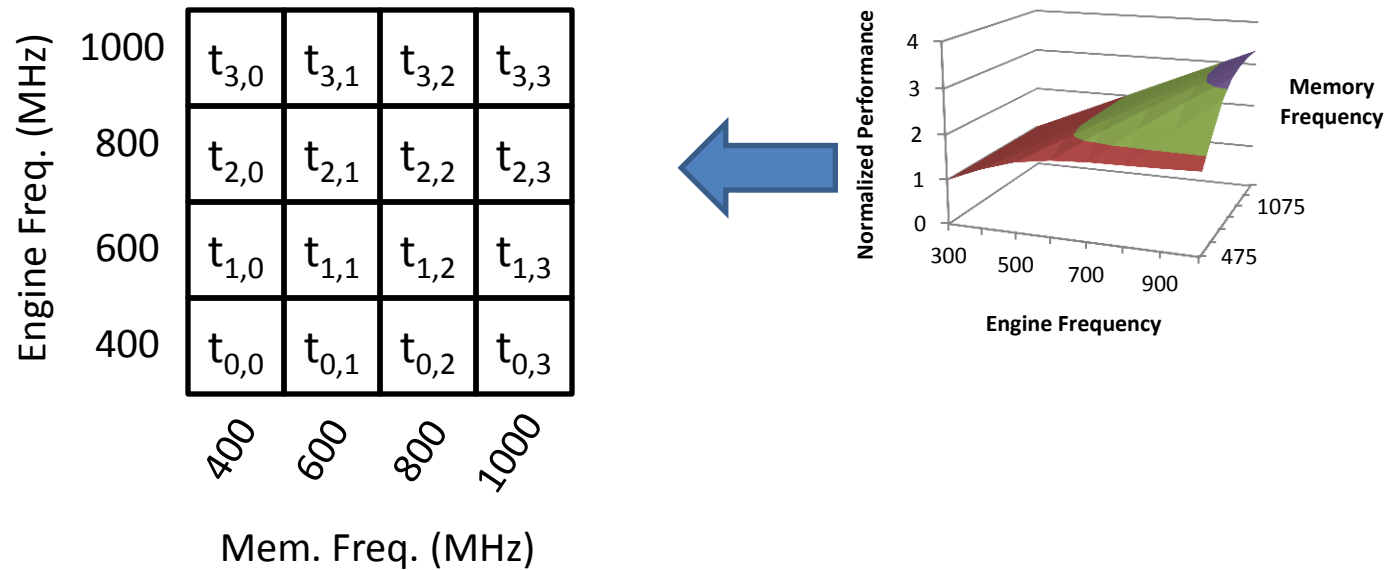
Backup Slides

Classifier: Neural Network Topology

- 3 layer, fully connected network
 - Input layer: linear
 - Number of neurons equals number of features
 - Hidden layer: sigmoid
 - Number of neurons equals number of clusters
 - Output layer: sigmoid
 - Number of neurons equals number of clusters



Execution Time Scaling Values



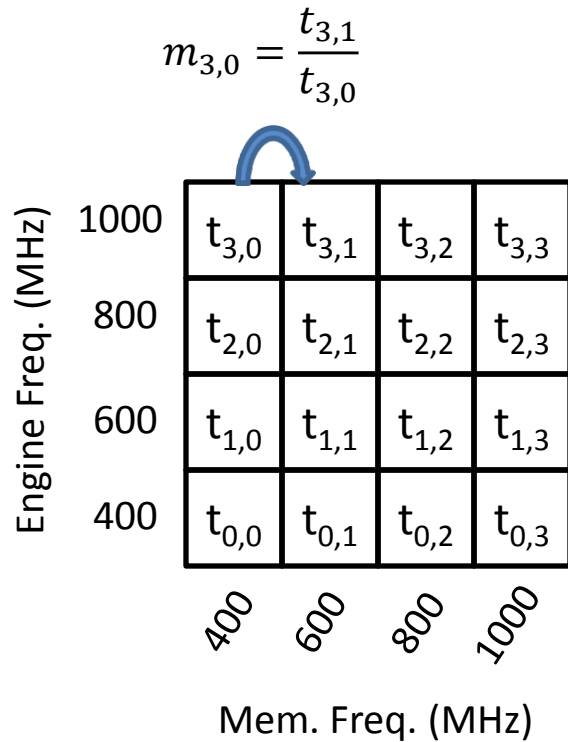
- Per kernel in training set
- Fixed CU count in this example

Execution Time Scaling Values

Engine Freq. (MHz)	1000	$t_{3,0}$	$t_{3,1}$	$t_{3,2}$	$t_{3,3}$
	800	$t_{2,0}$	$t_{2,1}$	$t_{2,2}$	$t_{2,3}$
	600	$t_{1,0}$	$t_{1,1}$	$t_{1,2}$	$t_{1,3}$
	400	$t_{0,0}$	$t_{0,1}$	$t_{0,2}$	$t_{0,3}$
		400	600	800	1000
		Mem. Freq. (MHz)			

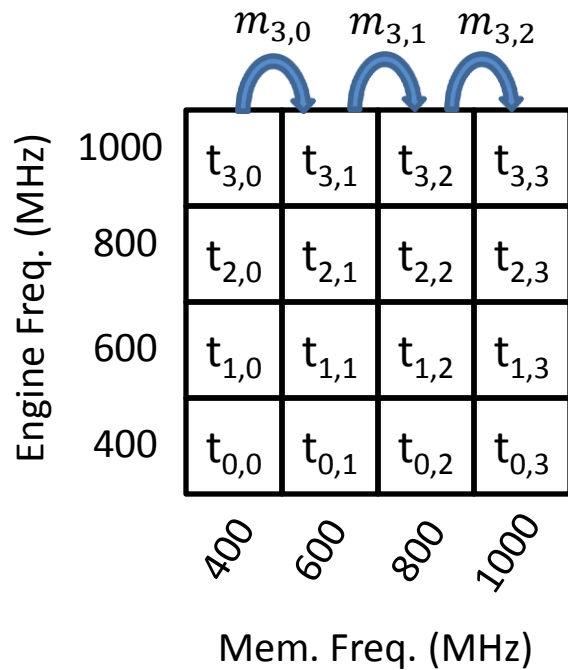
- Per kernel in training set
- Fixed CU count in this example

Execution Time Scaling Values



- Per kernel in training set
- Fixed CU count in this example

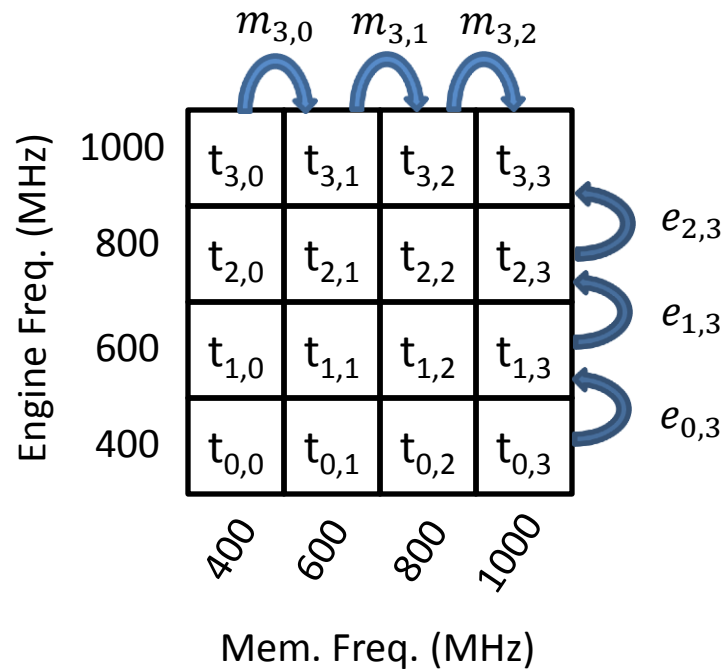
Execution Time Scaling Values



$$t_{3,3} = t_{3,0} \prod_{j=0}^2 m_{3,j}$$

- Per kernel in training set
- Fixed CU count in this example

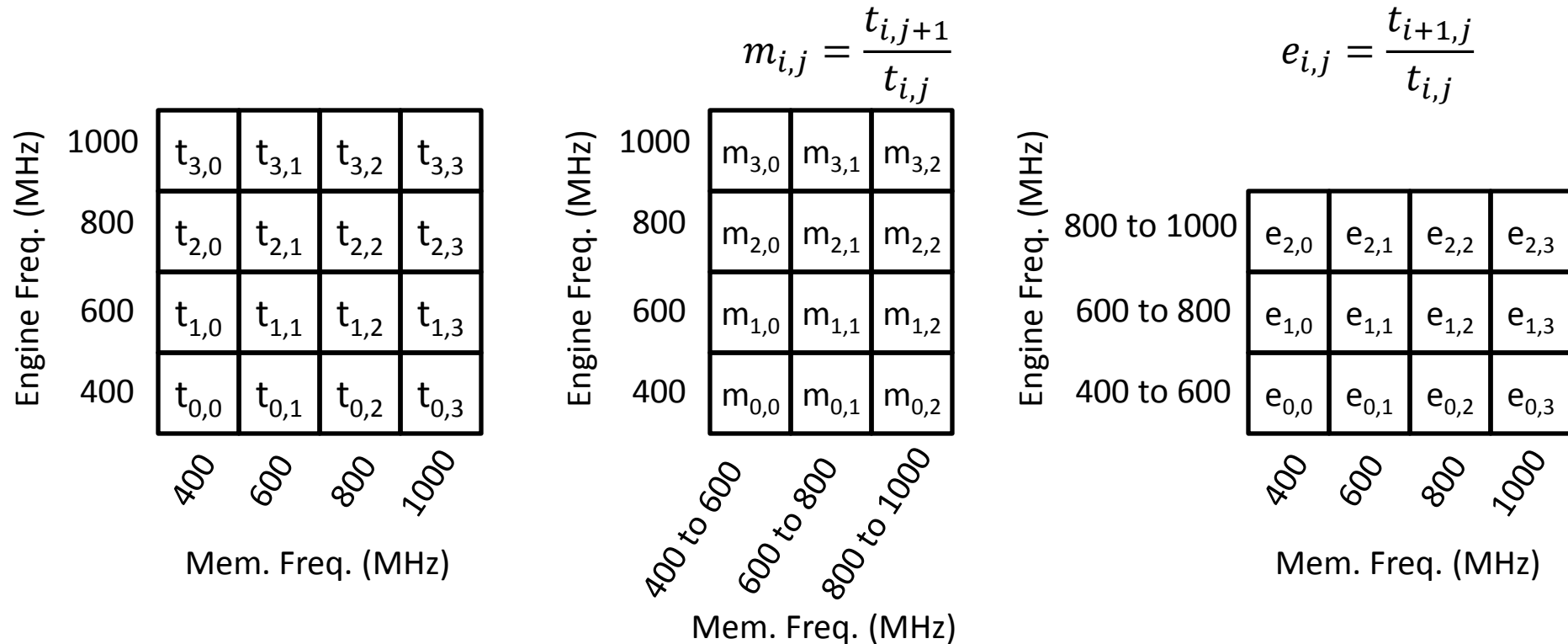
Execution Time Scaling Values



$$t_{3,3} = t_{3,0} \prod_{j=0}^2 m_{3,j}$$

- Per kernel in training set
- Fixed CU count in this example

Execution Time Scaling Values



- Per kernel in training set
- Fixed CU count in this example

K-means Clustering: the view from 10,000 feet

- Each kernel has a feature vector (x vector)

Kernel 1

$[x_0, x_1, x_2, \dots x_n]$

Kernel 2

$[x_0, x_1, x_2, \dots x_n]$

Kernel 3

$[x_0, x_1, x_2, \dots x_n]$

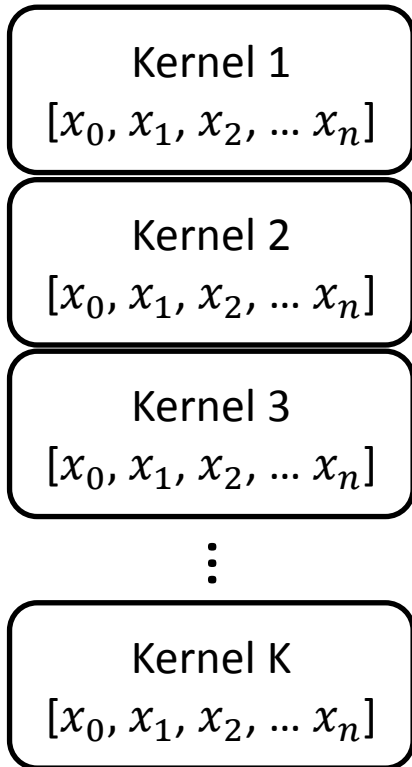
⋮

Kernel K

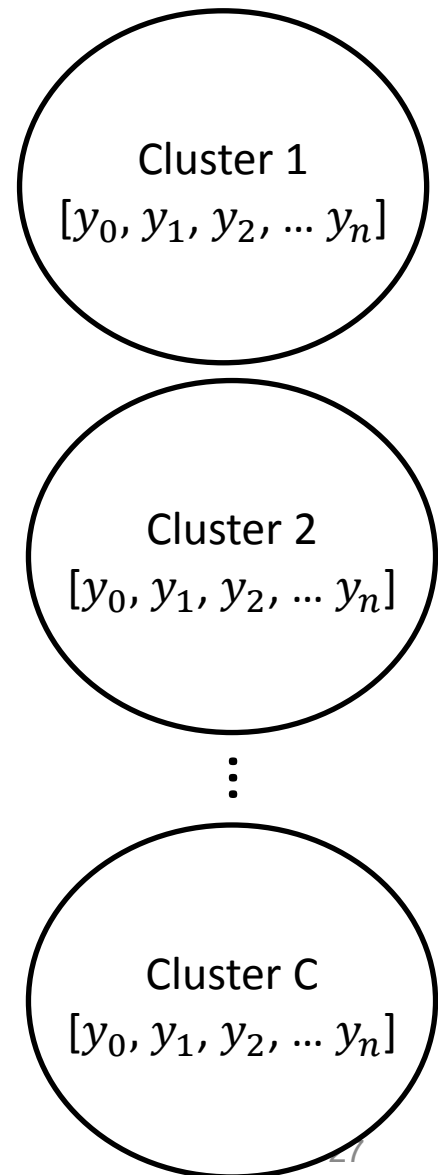
$[x_0, x_1, x_2, \dots x_n]$

K-means Clustering: the view from 10,000 feet

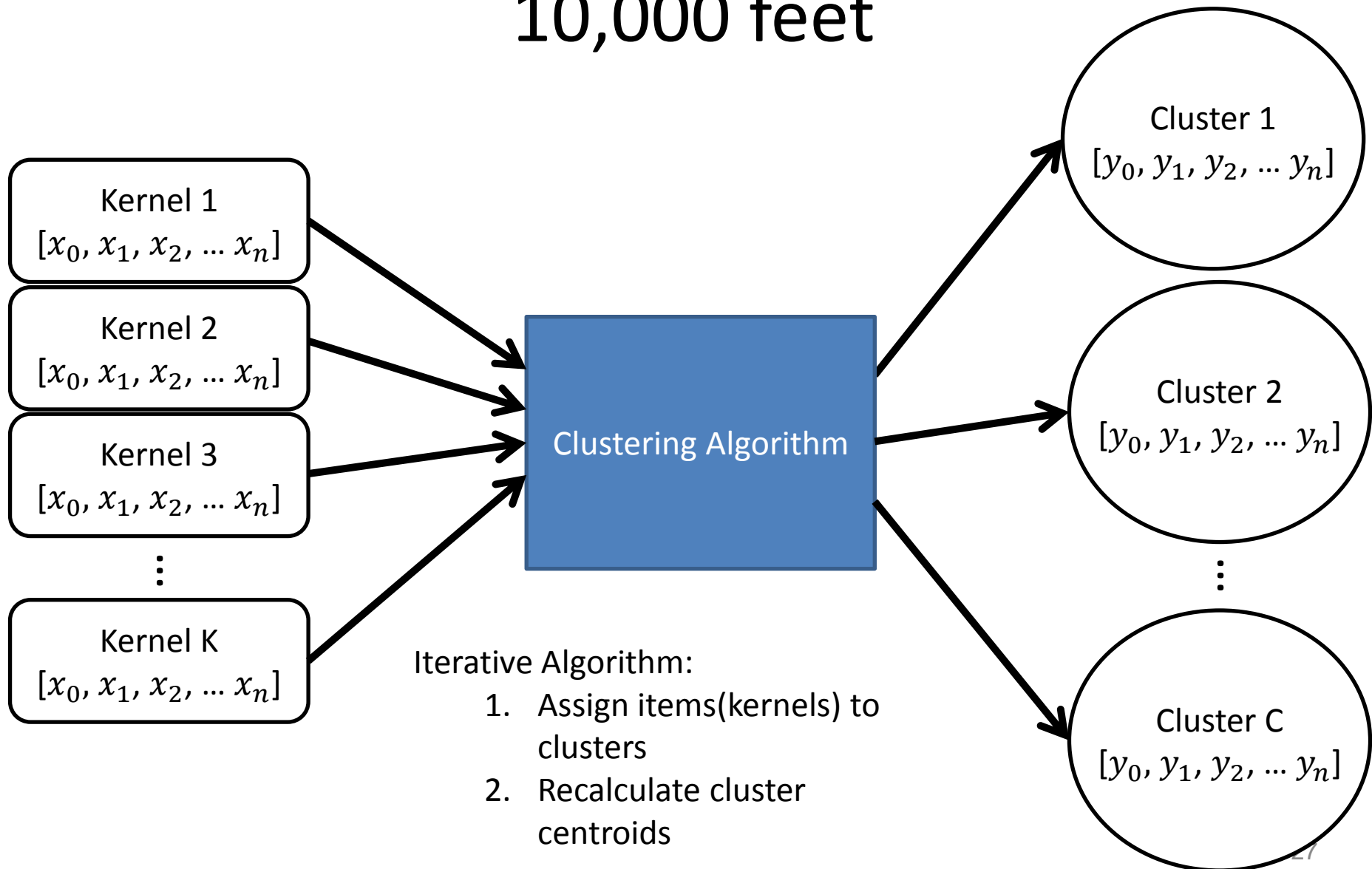
- Each kernel has a feature vector (x vector)



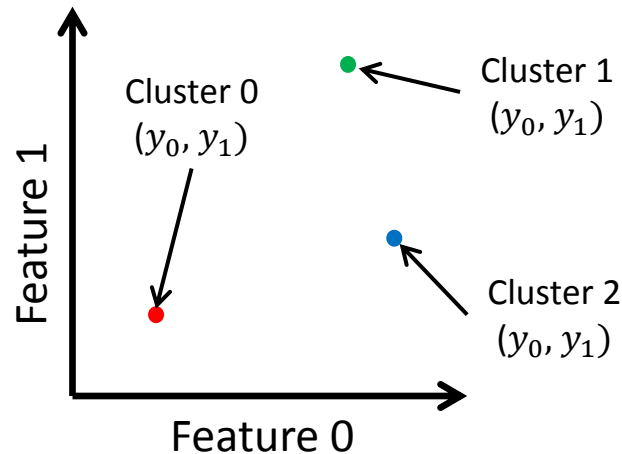
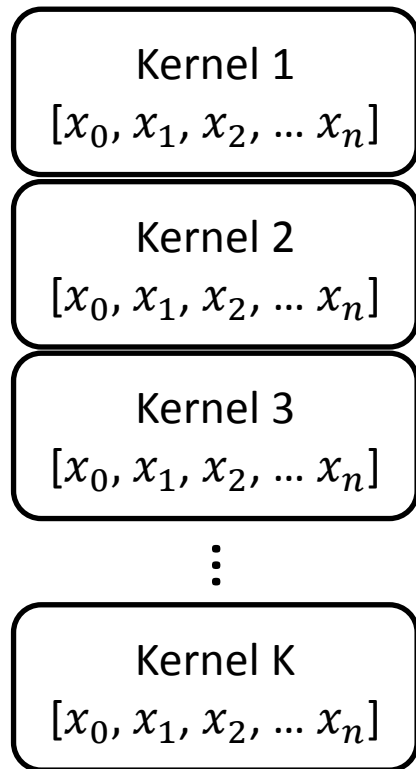
- Each kernel has a cluster has a centroid (y vector)



K-means Clustering: the view from 10,000 feet

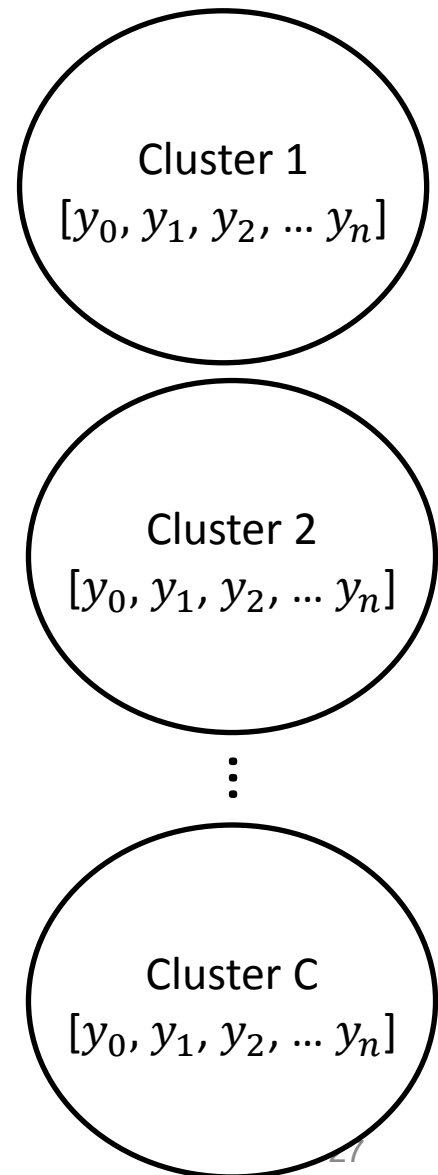


K-means Clustering: the view from 10,000 feet

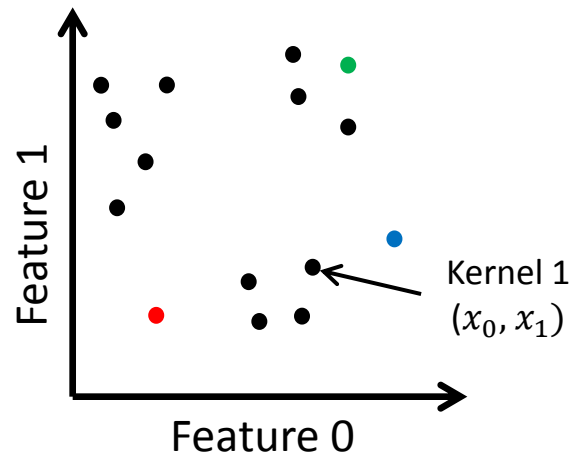
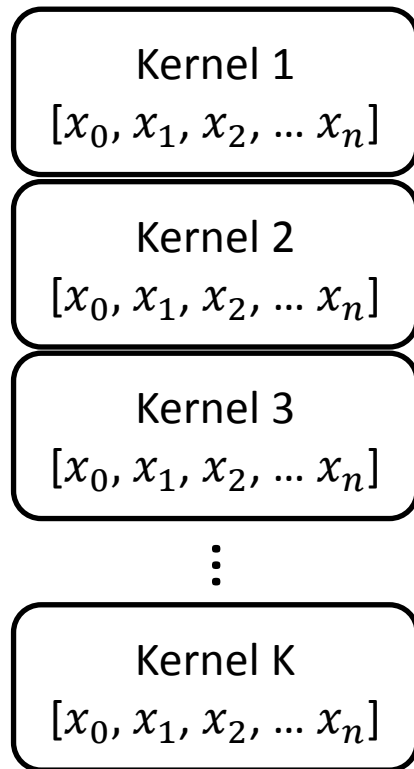


Iterative Algorithm:

1. Assign items(kernels) to clusters
2. Recalculate cluster centroids

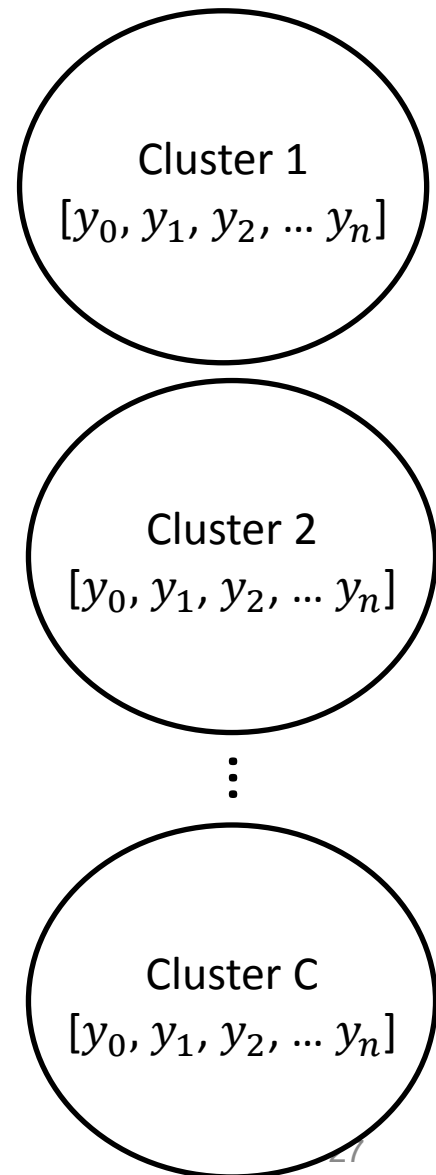


K-means Clustering: the view from 10,000 feet

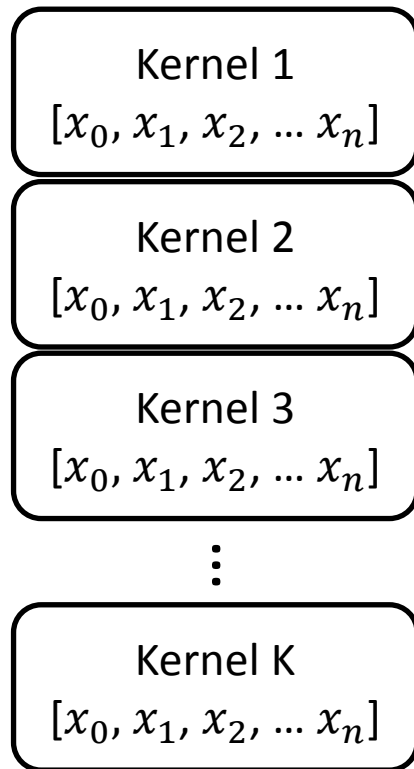


Iterative Algorithm:

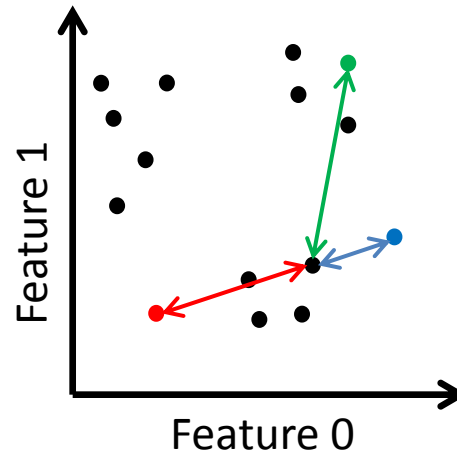
1. Assign items(kernels) to clusters
2. Recalculate cluster centroids



K-means Clustering: the view from 10,000 feet

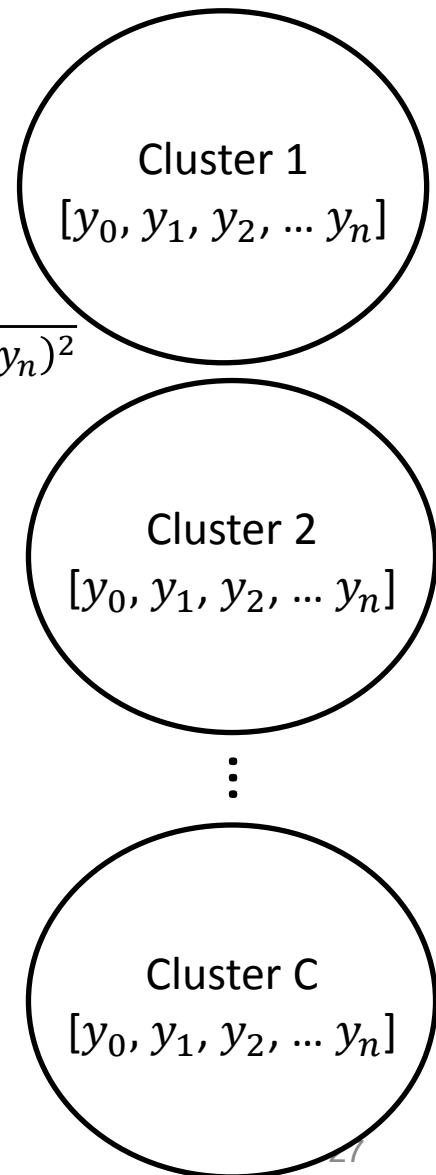


$$\text{distance} = \sqrt{(x_0 - y_0)^2 + (x_1 - y_1)^2 + \dots + (x_n - y_n)^2}$$

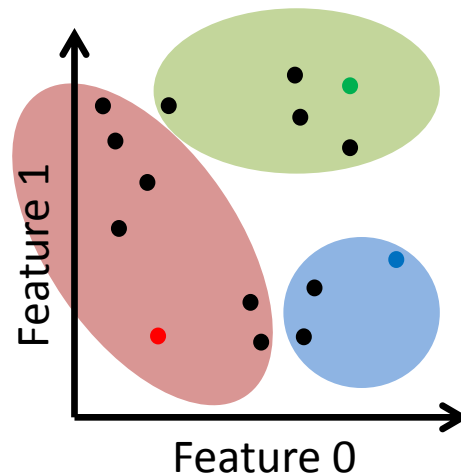
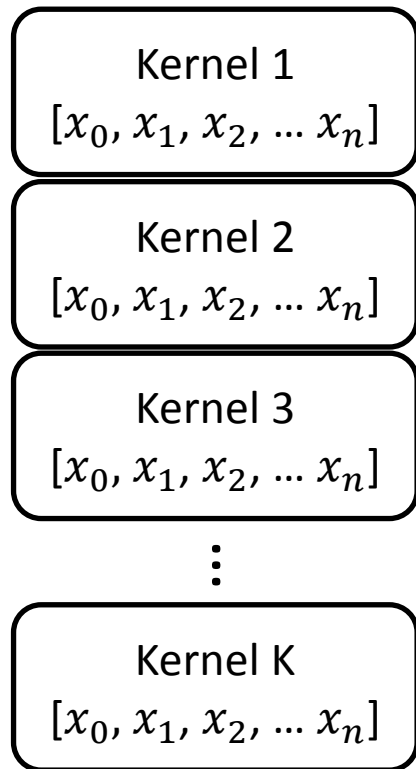


Iterative Algorithm:

1. Assign items(kernels) to clusters
2. Recalculate cluster centroids

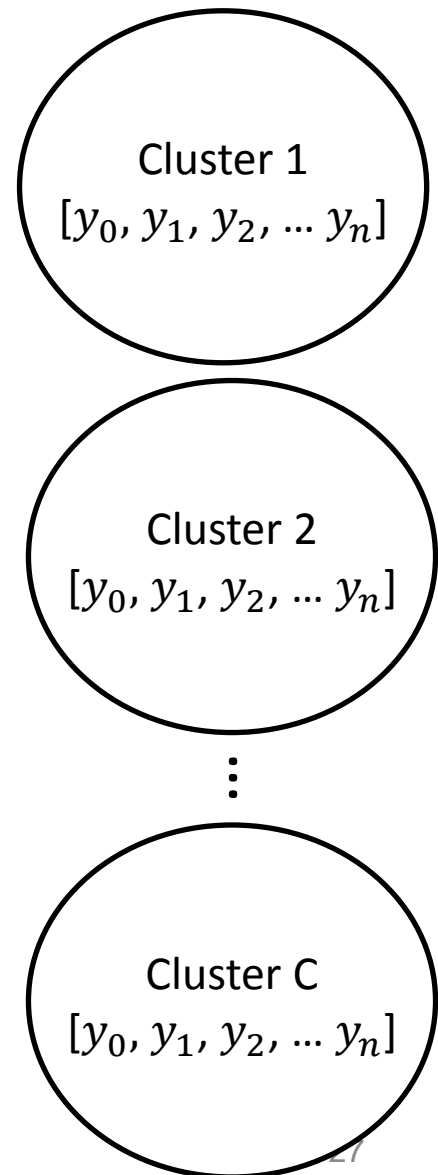


K-means Clustering: the view from 10,000 feet

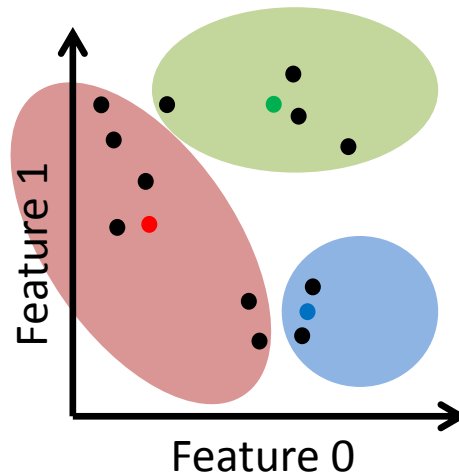
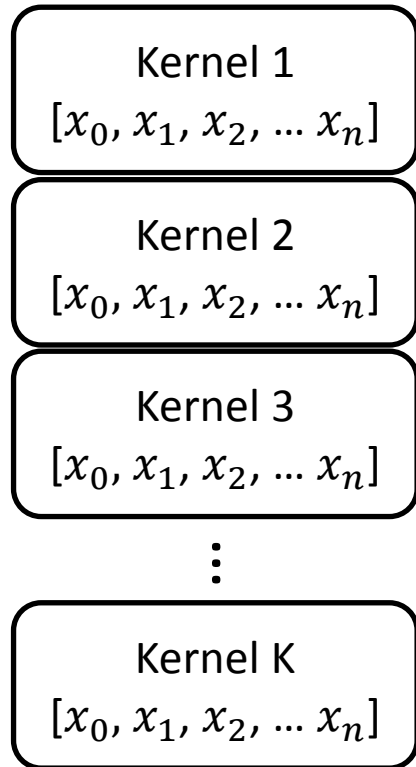


Iterative Algorithm:

1. Assign items(kernels) to clusters
2. Recalculate cluster centroids

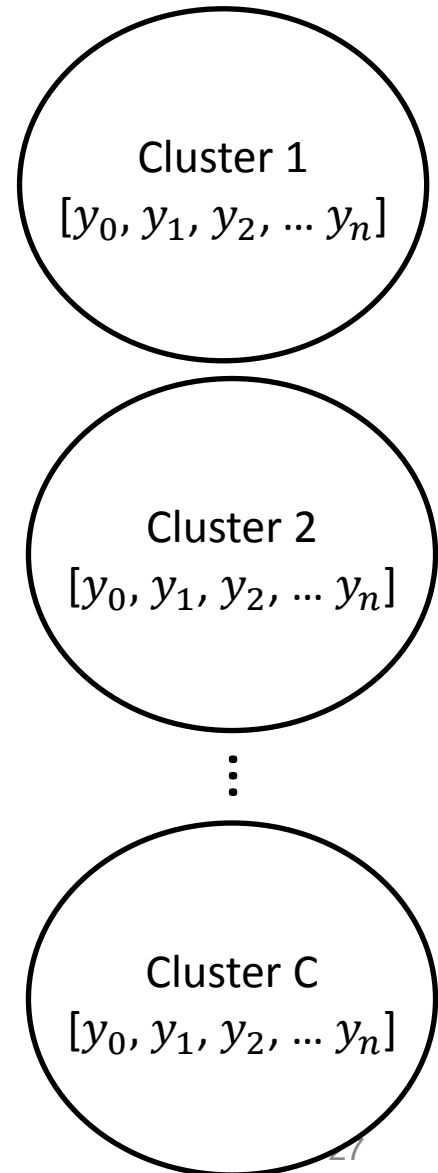


K-means Clustering: the view from 10,000 feet

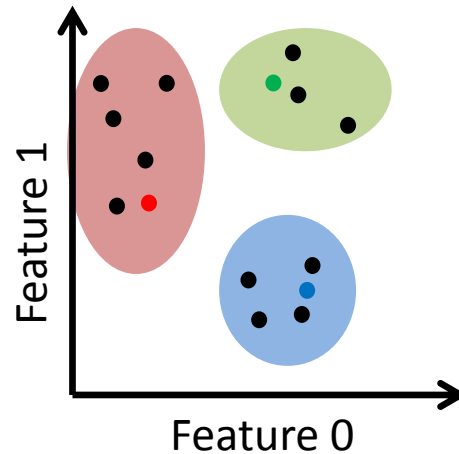
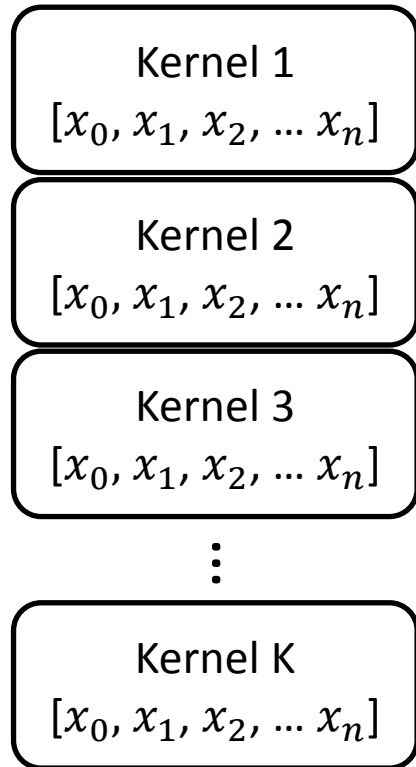


Iterative Algorithm:

1. Assign items(kernels) to clusters
2. Recalculate cluster centroids

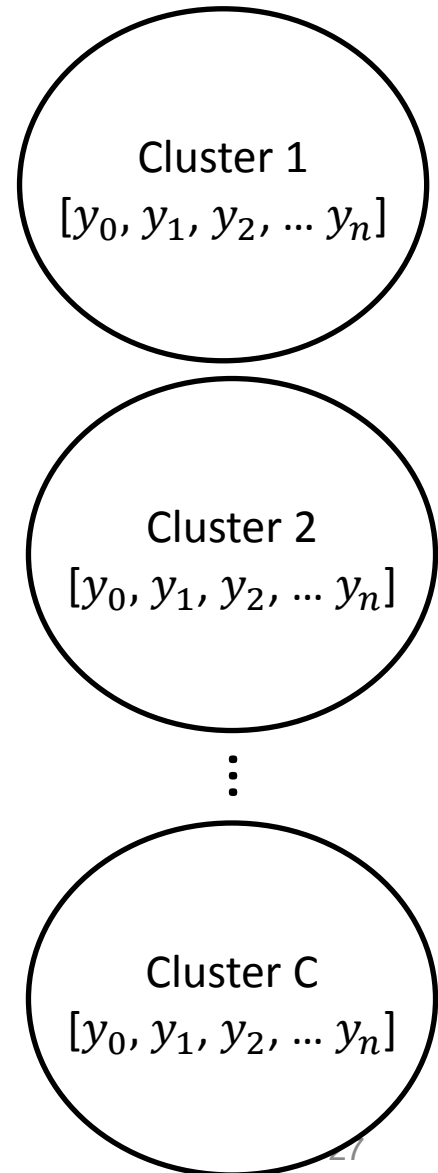


K-means Clustering: the view from 10,000 feet

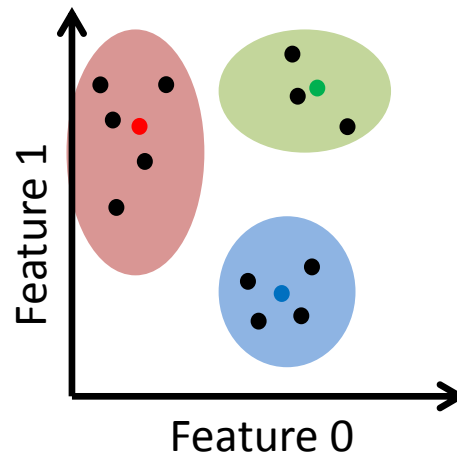
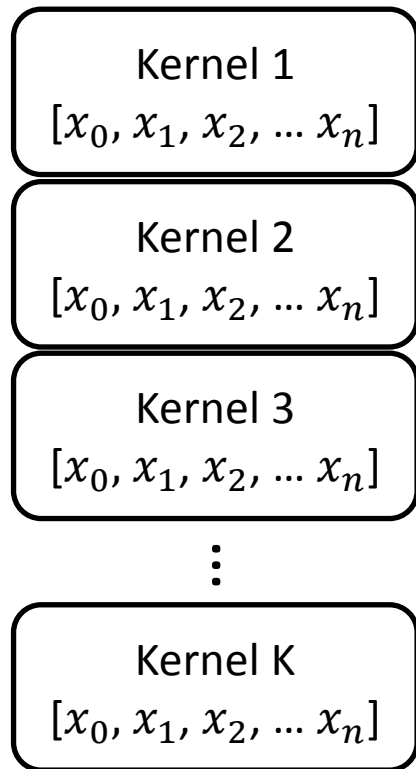


Iterative Algorithm:

1. Assign items(kernels) to clusters
2. Recalculate cluster centroids

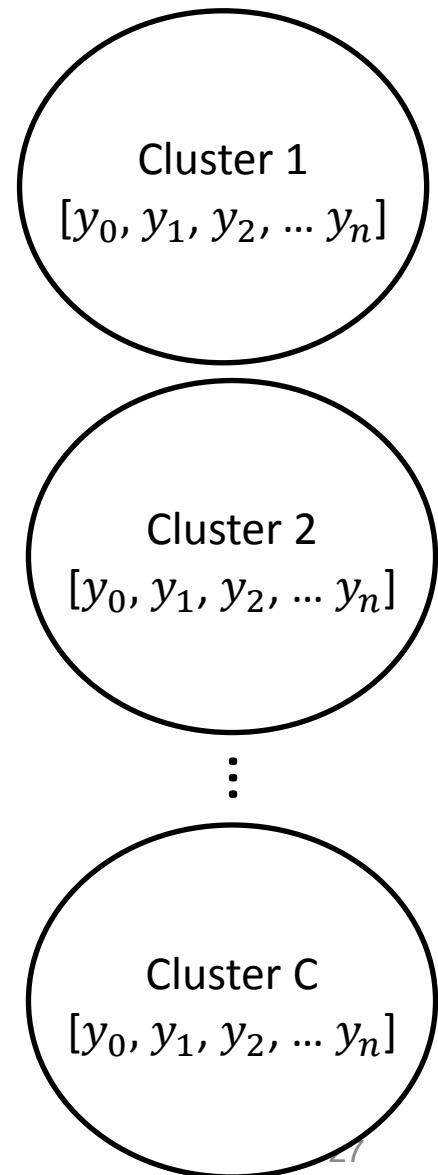


K-means Clustering: the view from 10,000 feet

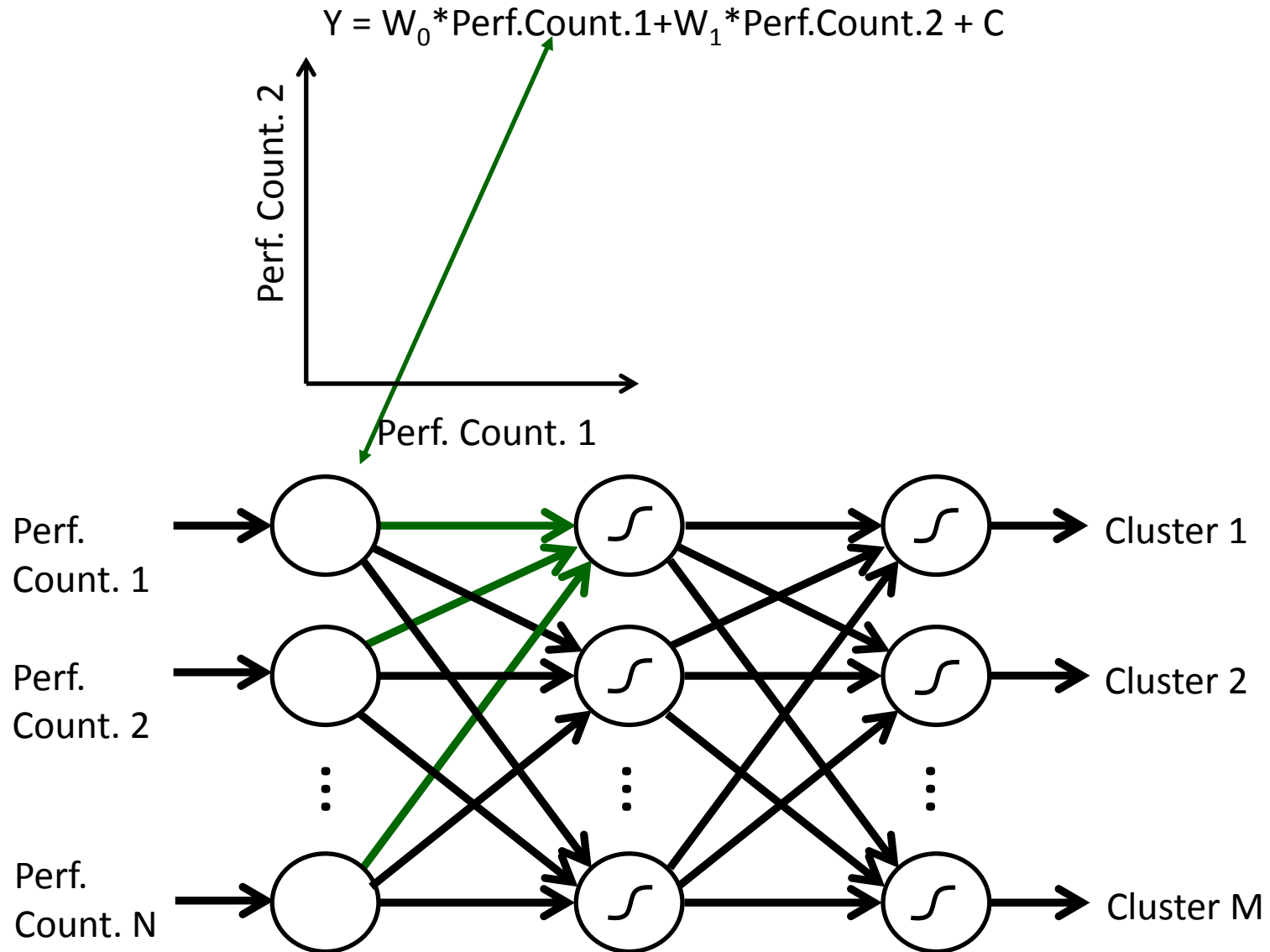


Iterative Algorithm:

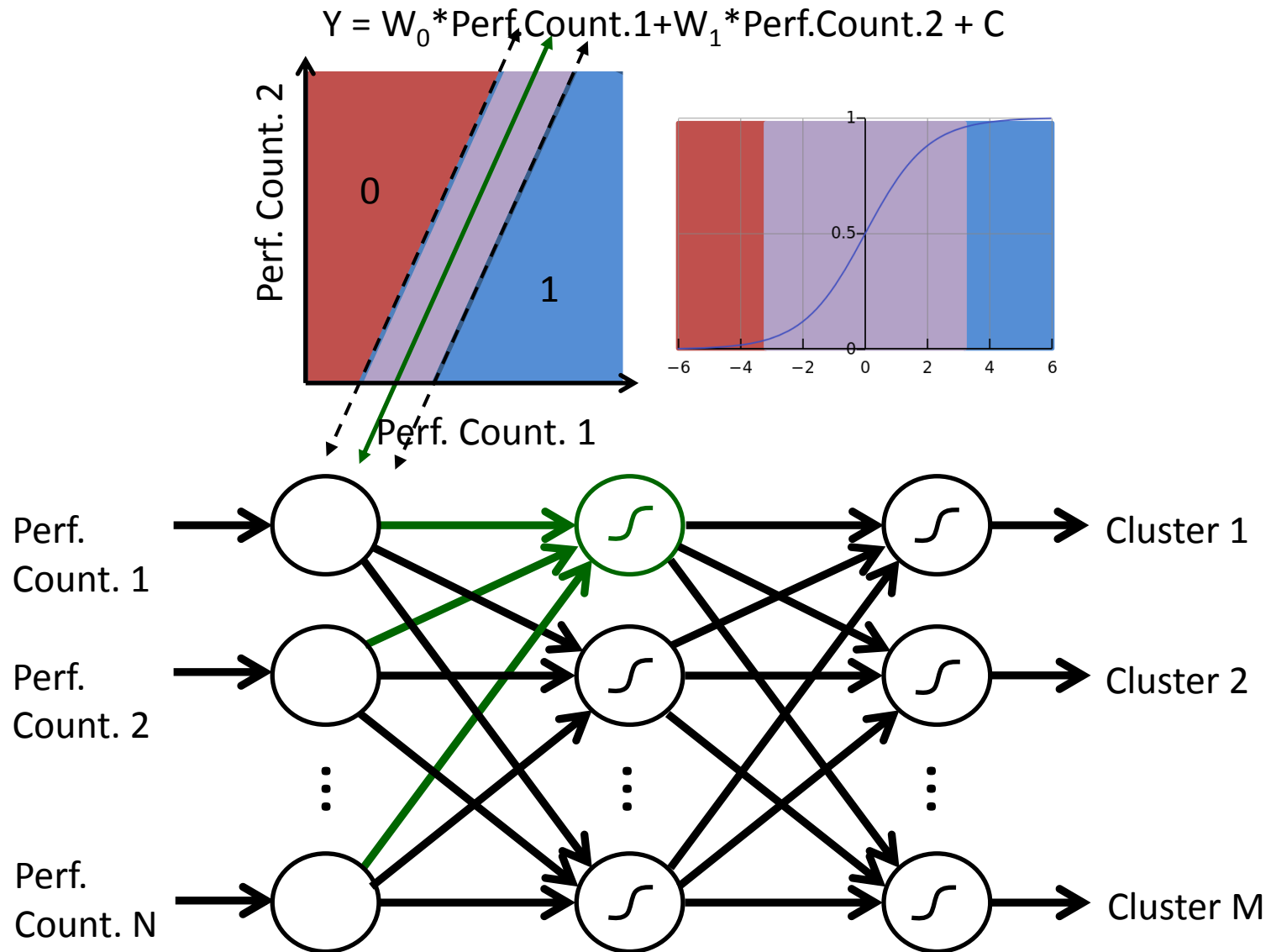
1. Assign items(kernels) to clusters
2. Recalculate cluster centroids



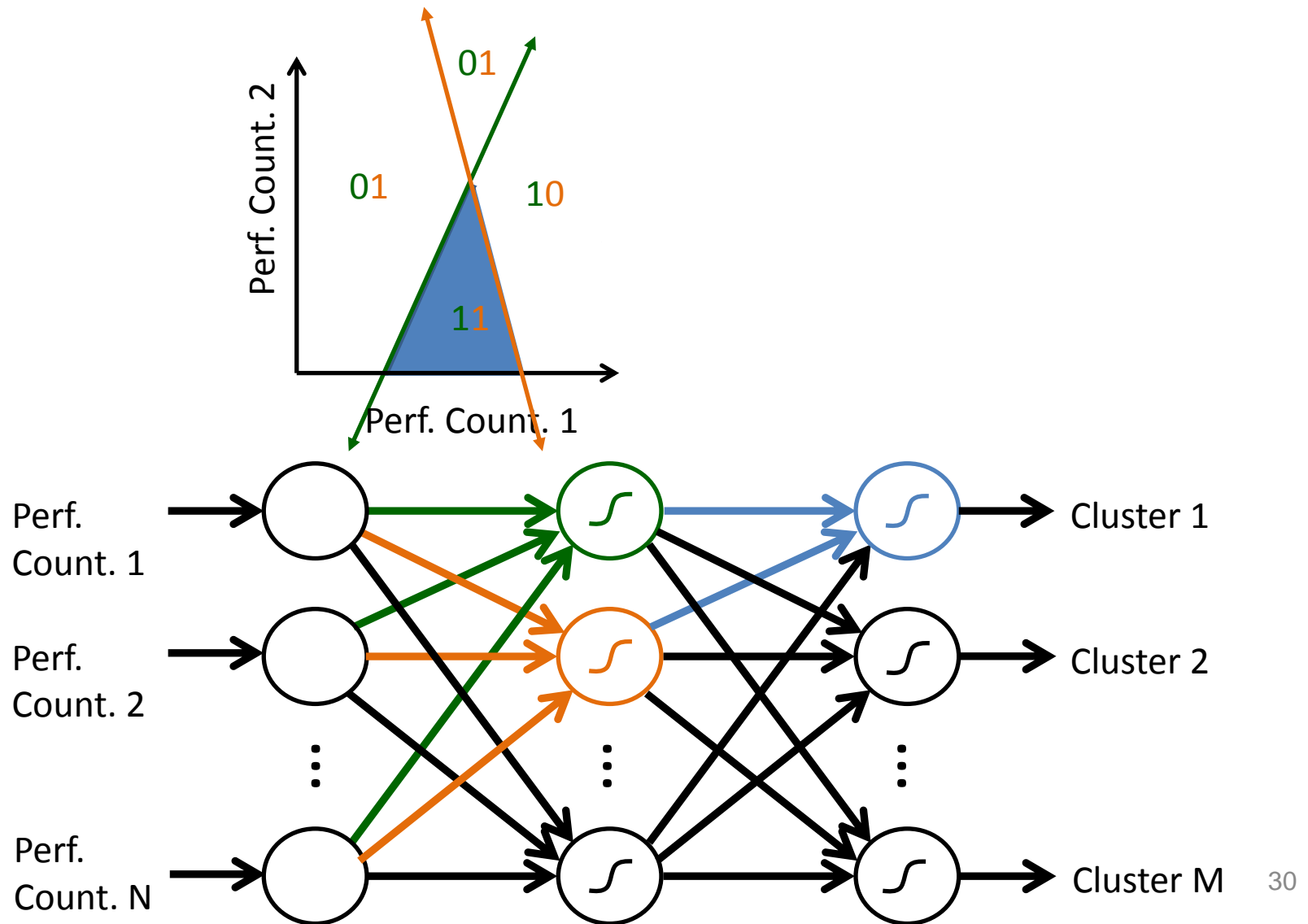
Neural Network



Neural Network



Neural Network



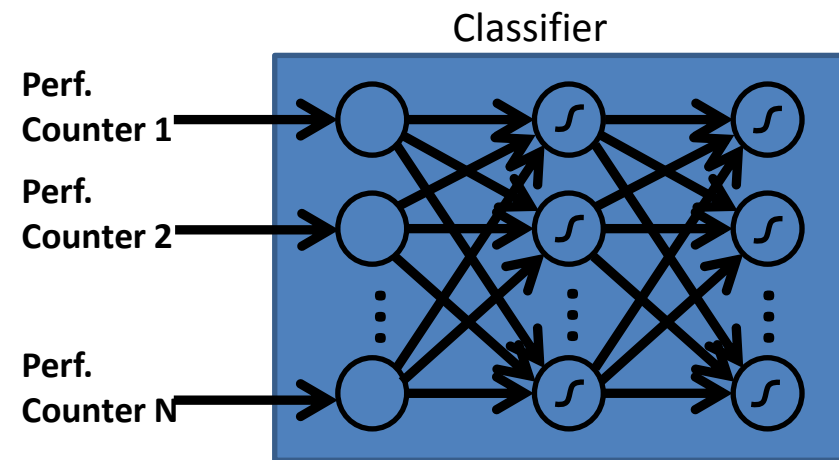
Putting It All Together

**Perf.
Counter 1**

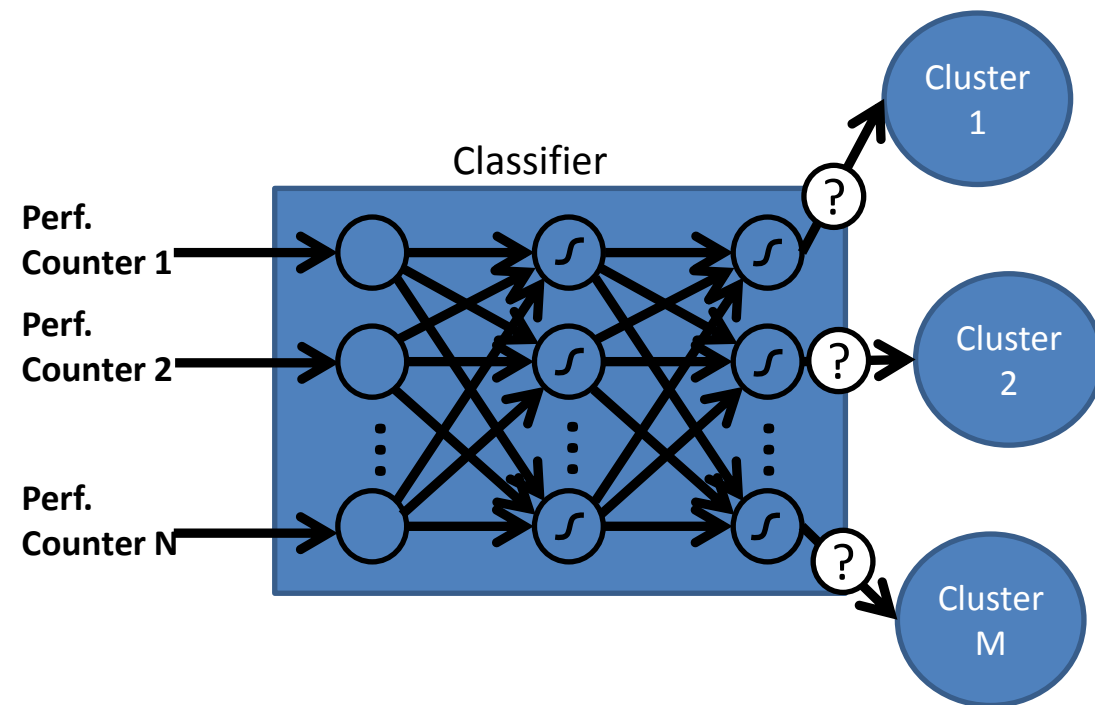
**Perf.
Counter 2**

**Perf.
Counter N**

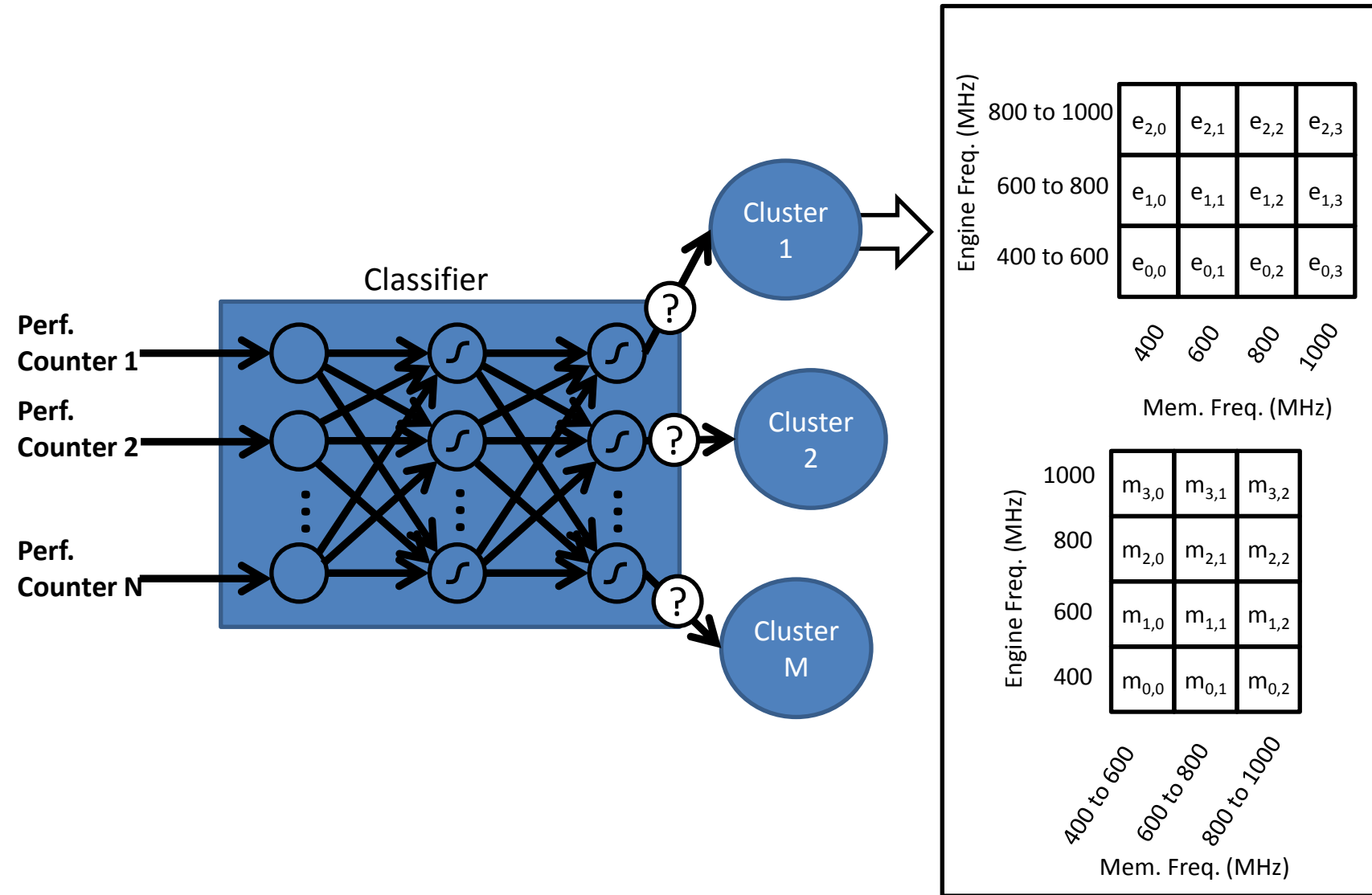
Putting It All Together



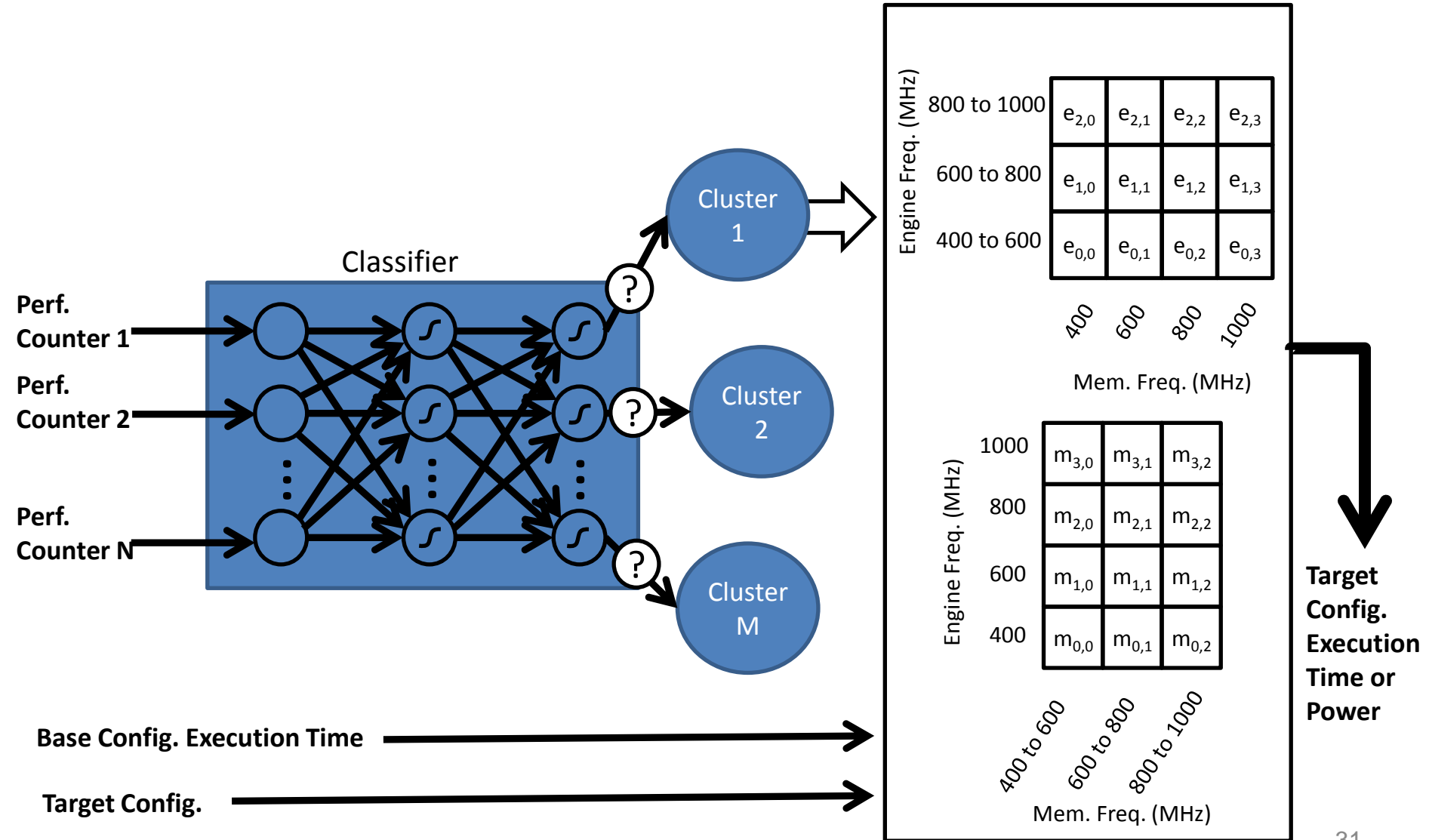
Putting It All Together



Putting It All Together



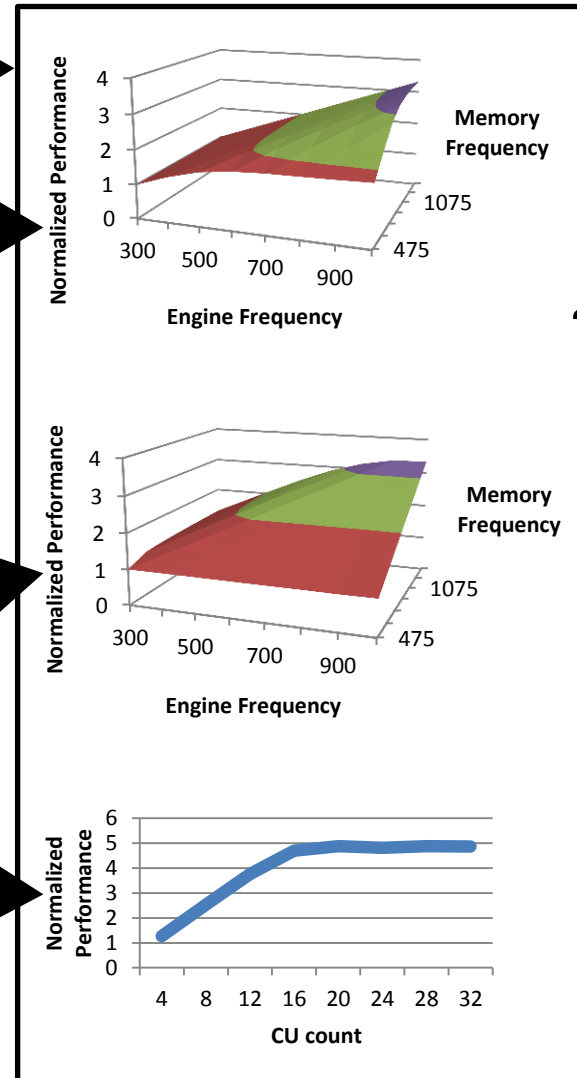
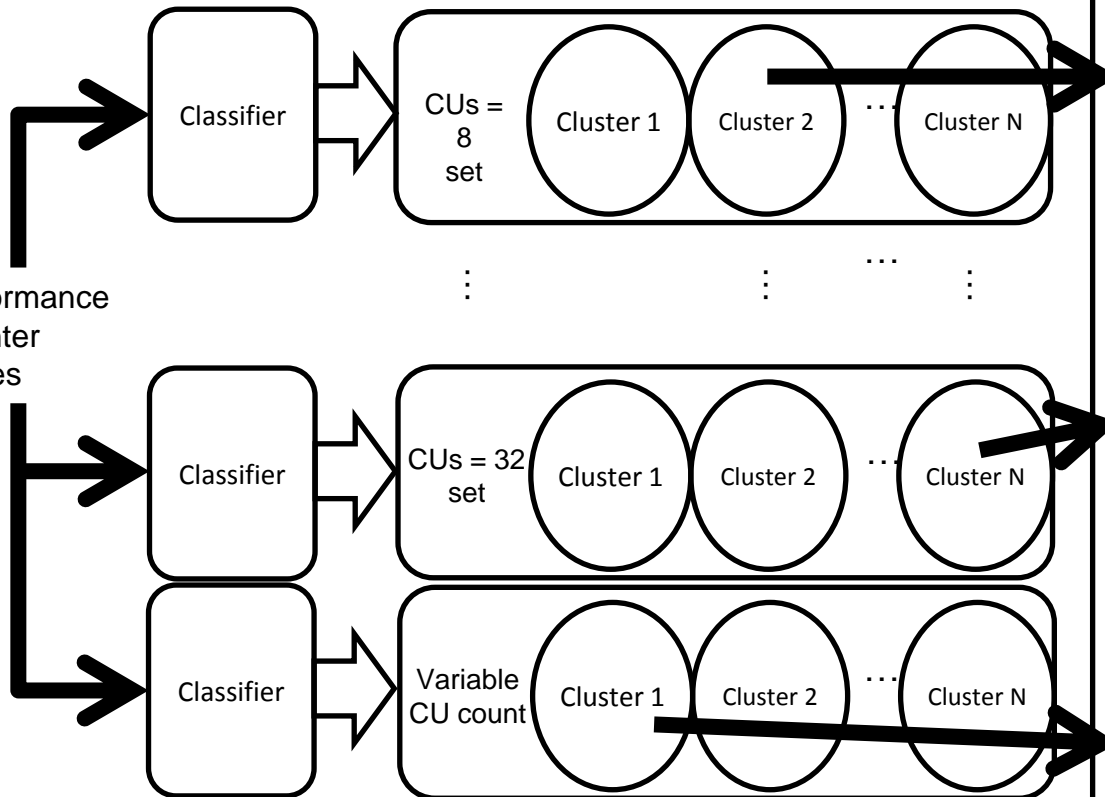
Putting It All Together



Model Architecture

Base Config.
Exec. Time &
Target Config.

Performance
Counter
Values



Target
Config.
Exec.
Time