

# Ptolemy: Architecture Support for Robust Deep Learning

**Yiming Gan**

Department of Computer Science,  
University of Rochester

with

Yuxian Qiu, Shanghai Jiao Tong University  
Jingwen Leng, Shanghai Jiao Tong University  
Minyi Guo, Shanghai Jiao Tong University  
Yuhao Zhu, University of Rochester

<https://github.com/Ptolemy-dl/Ptolemy>



# Deep Learning: Not Robust



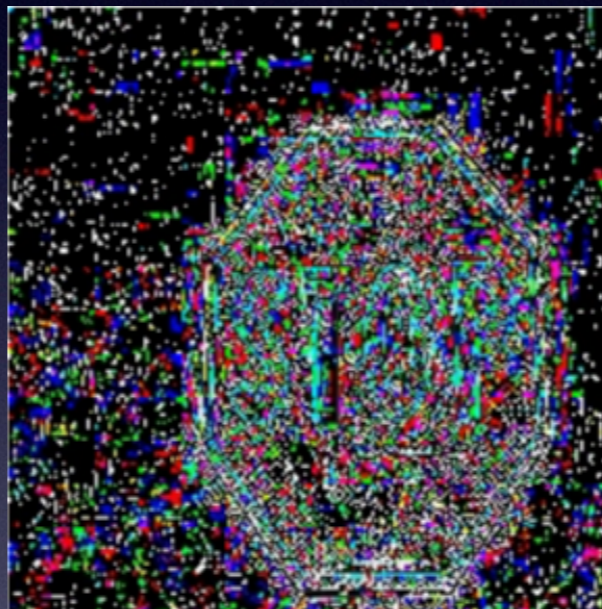
# Deep Learning: Not Robust

Legitimate Example



+

Perturbation



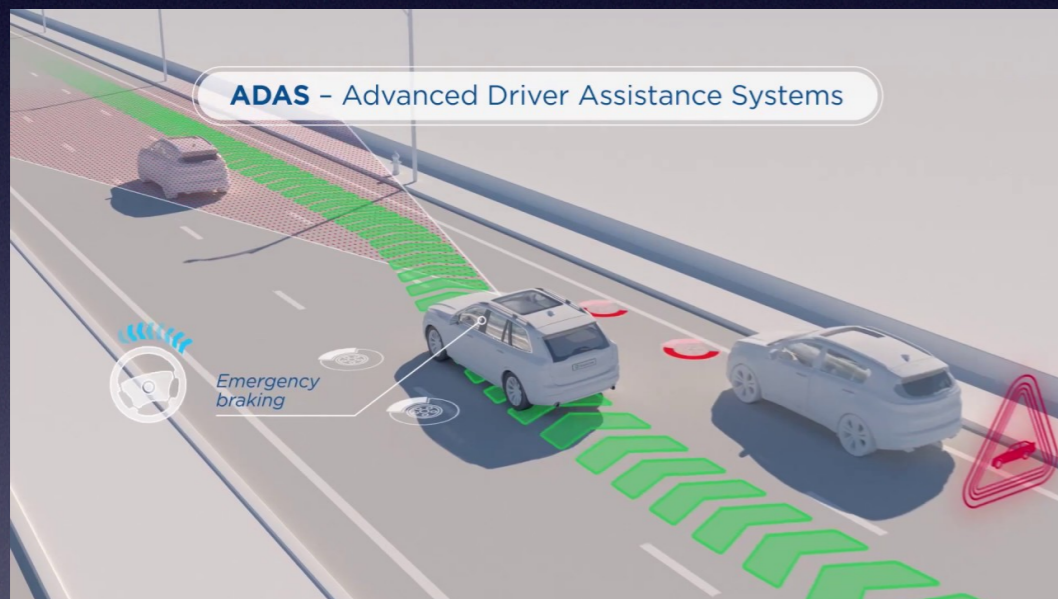
=

Adversarial Example

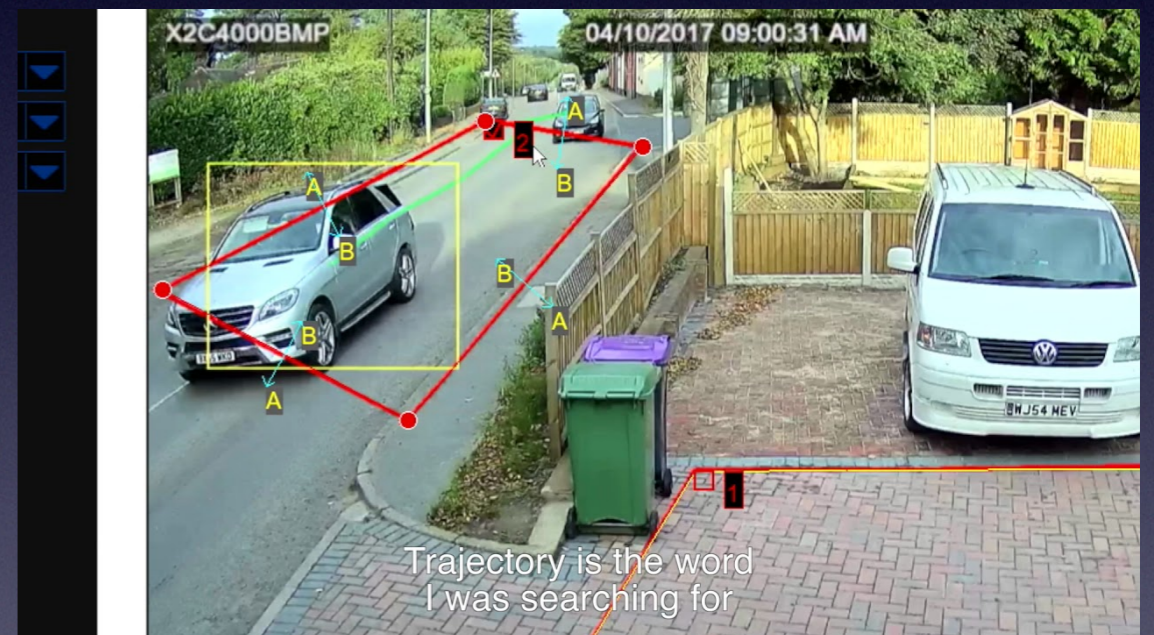


# Mission Critical System

## ADAS



## Security Cameras



# Robust Deep Learning Requirements

- Accurately detect adversarial examples

# Robust Deep Learning Requirements

- Accurately detect adversarial examples

+

- Do not bring large overhead on system performance

# Robust Deep Learning Requirements

- Accurately detect adversarial examples

+

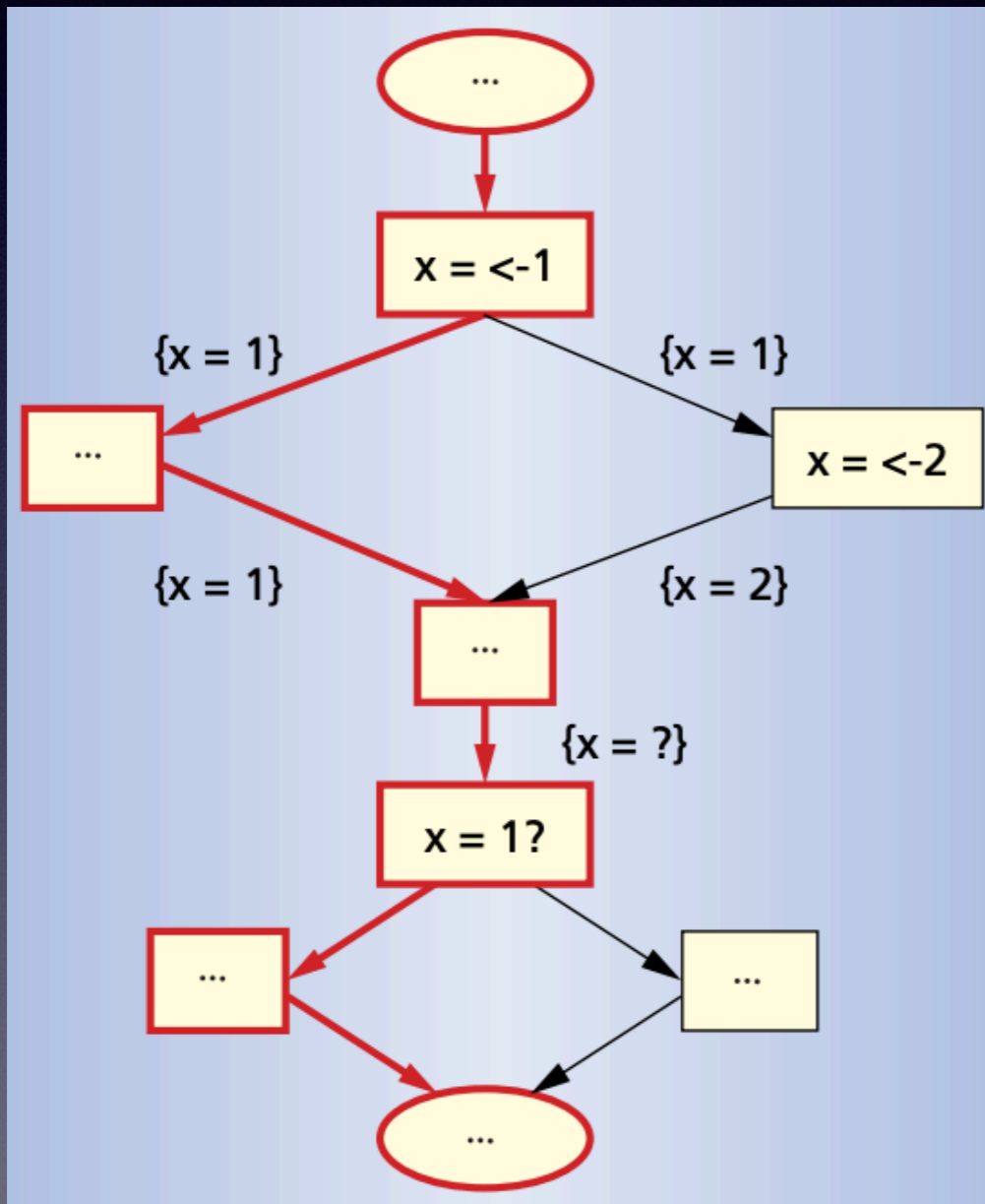
- Do not bring large overhead on system performance

=

**Ptolemy**

# Hot Path

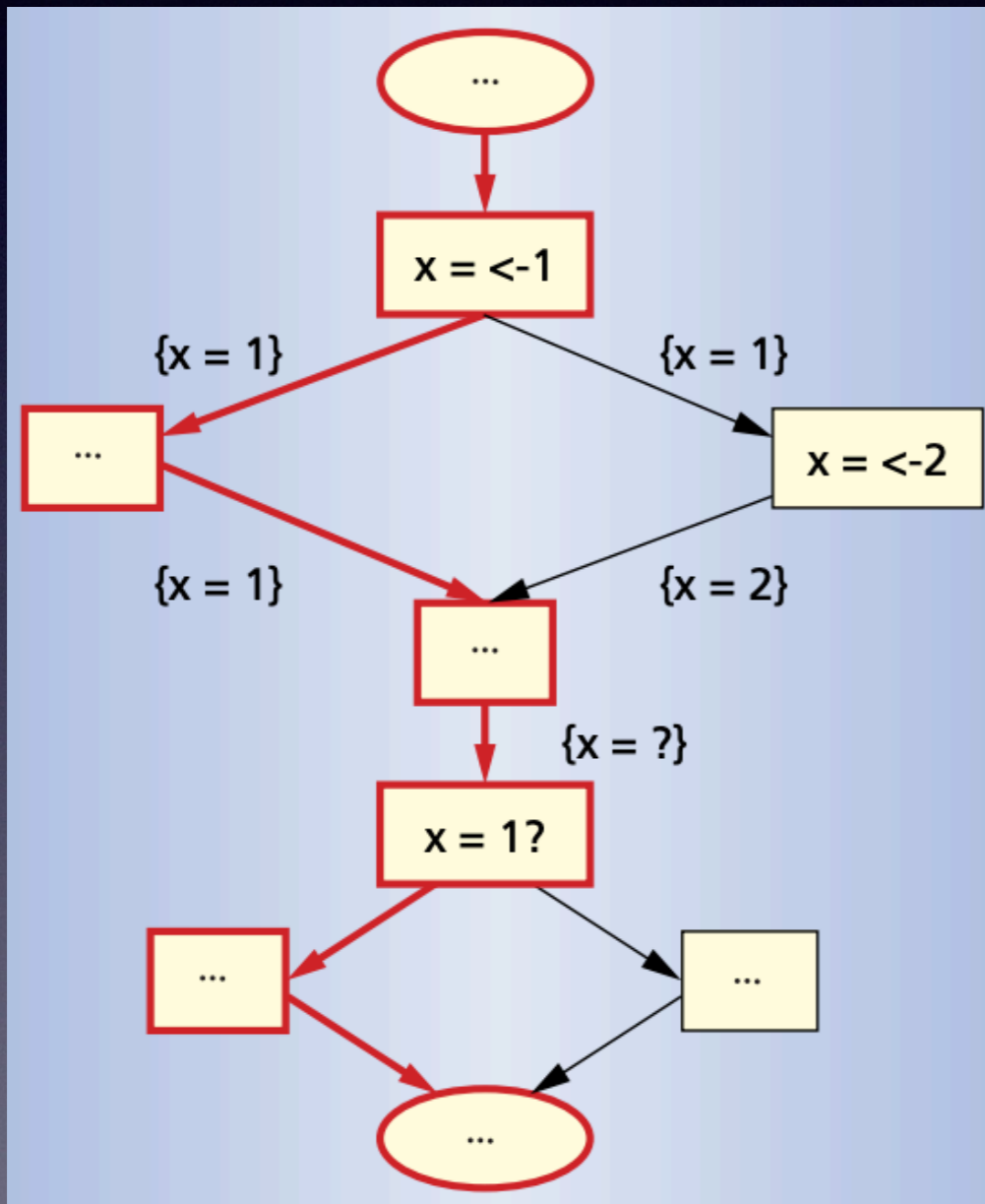
## Traditional Software





# Hot Path

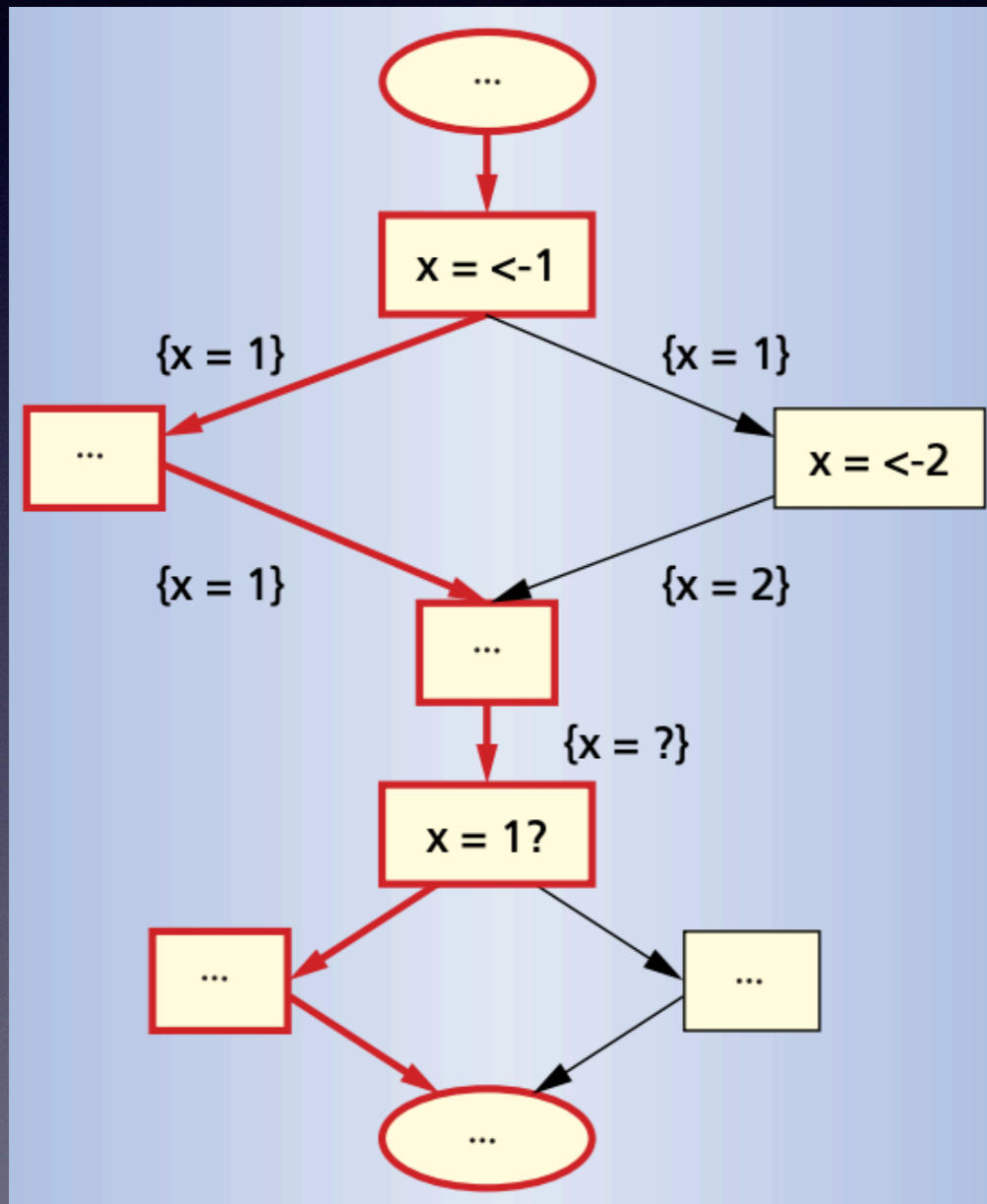
## Traditional Software



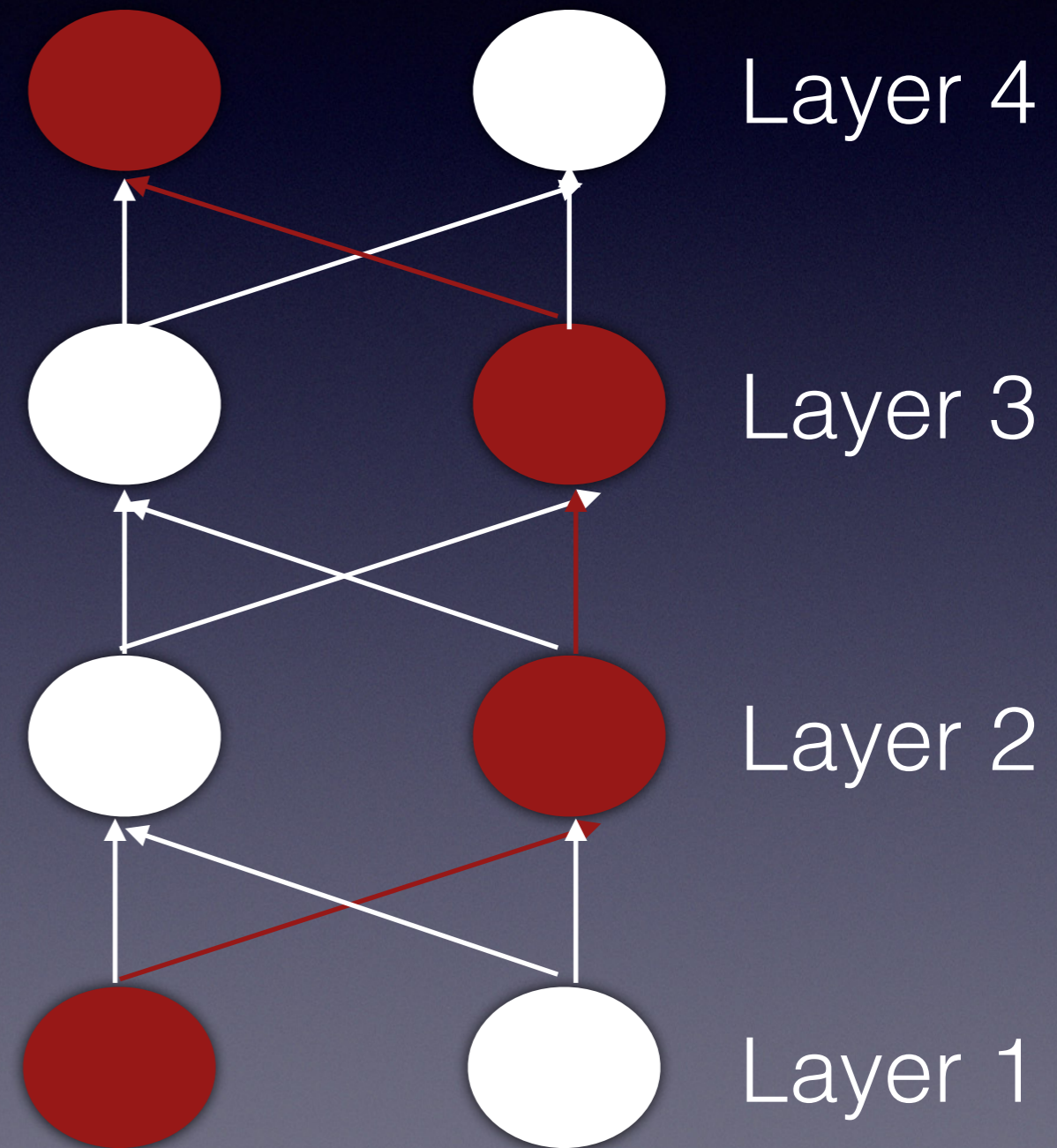
- Measure Program Behavior
- Optimizing Program
- Debugging

# Hot Path

## Traditional Software



## Deep Learning



# Defining Important Neuron

Input  
Feature Map

Weights

Output  
Feature Map

0.3	0.4	0.2	1.0	0.1
-----	-----	-----	-----	-----

**x**

0.2	0.2	0.3
0.3	0.2	0.4
0.4	0.1	0.2
-0.1	0.09	0.1
-1.0	2.1	0.5

**=**

0.06
0.46
0.44

# Defining Important Neuron

Input  
Feature Map

Weights

Output  
Feature Map

0.3	0.4	0.2	1.0	0.1
-----	-----	-----	-----	-----

**x**

0.2	0.2	0.3
0.3	0.2	0.4
0.4	0.1	0.2
-0.1	0.09	0.1
-1.0	2.1	0.5

**=**

0.06
<b>0.46</b>
0.44

# Defining Important Neuron

Input  
Feature Map

0.3	0.4	0.2	1.0	0.1
-----	-----	-----	-----	-----

**x**

Weights

0.2	0.2	0.3
0.3	0.2	0.4
0.4	0.1	0.2
-0.1	0.09	0.1
-1.0	2.1	0.5

Output  
Feature Map

=

0.06
<b>0.46</b>
0.44

# Defining Important Neuron

Input  
Feature Map

Weights

Output  
Feature Map

0.3	0.4	0.2	<b>1.0</b>	<b>0.1</b>
-----	-----	-----	------------	------------

**x**

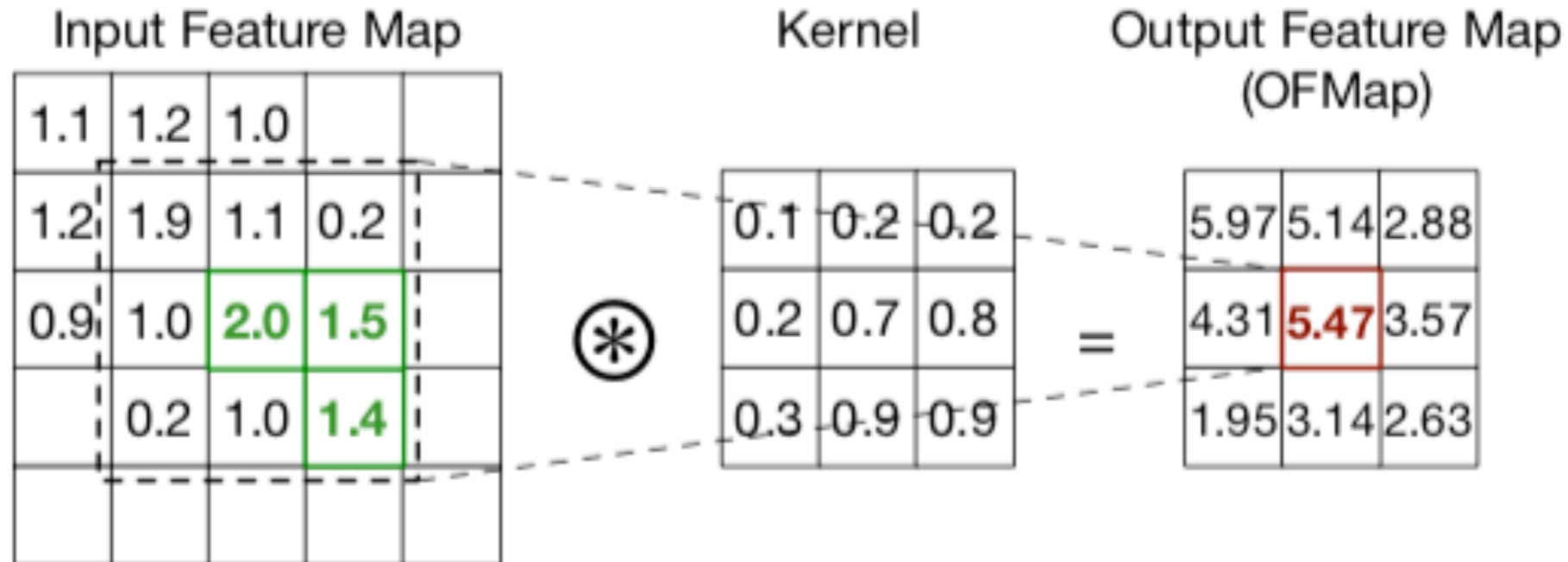
0.2	0.2	0.3
0.3	0.2	0.4
0.4	0.1	0.2
-0.1	<b>0.09</b>	0.1
-1.0	<b>2.1</b>	0.5

=

0.06
<b>0.46</b>
0.44

# Defining Important Neuron

## Important Neuron Extraction in Convolution Layer



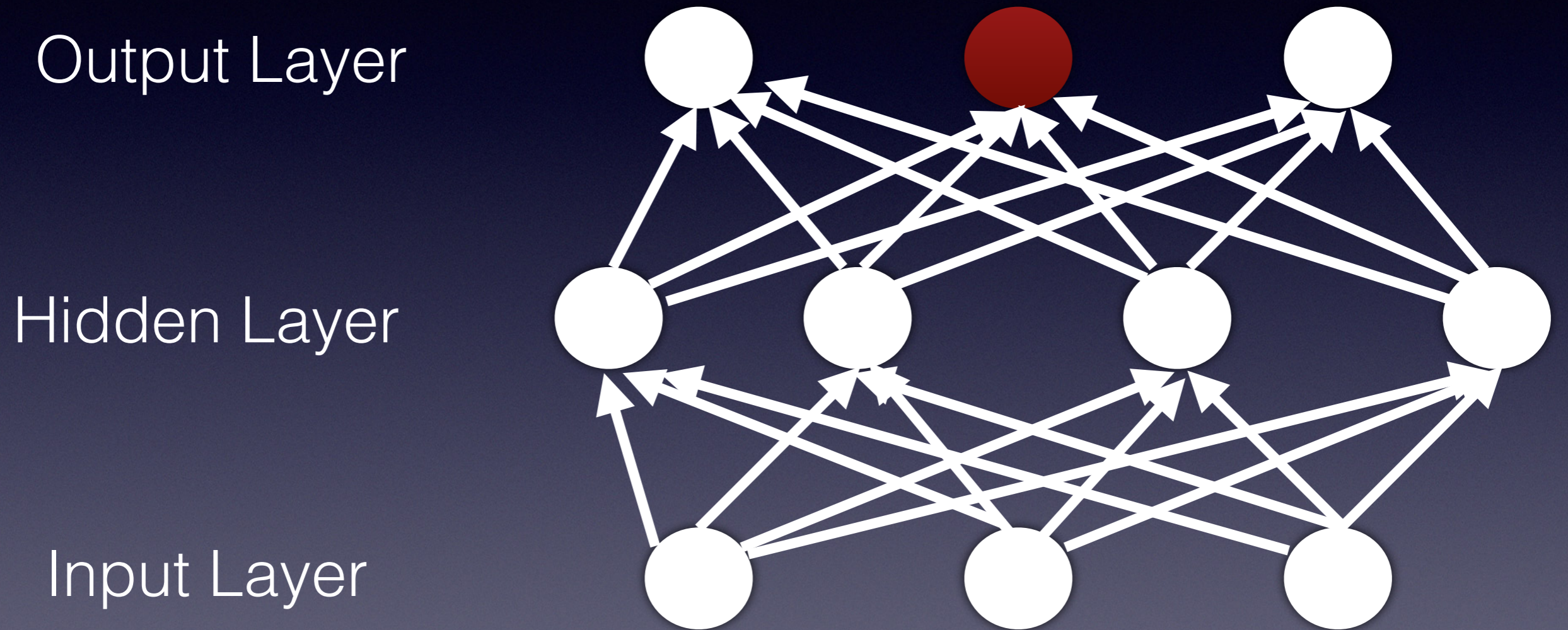
$$5.47 = 2.0 \times 0.7 + 1.4 \times 0.9 + 1.5 \times 0.8 + 1.0 \times 0.9 + \dots$$

$$2.0 \times 0.7 + 1.4 \times 0.9 + 1.5 \times 0.8 > 0.6 \times 5.47, \text{ assuming } \theta = 0.6$$

Important Neurons in the OFMap (identified before): 5.47

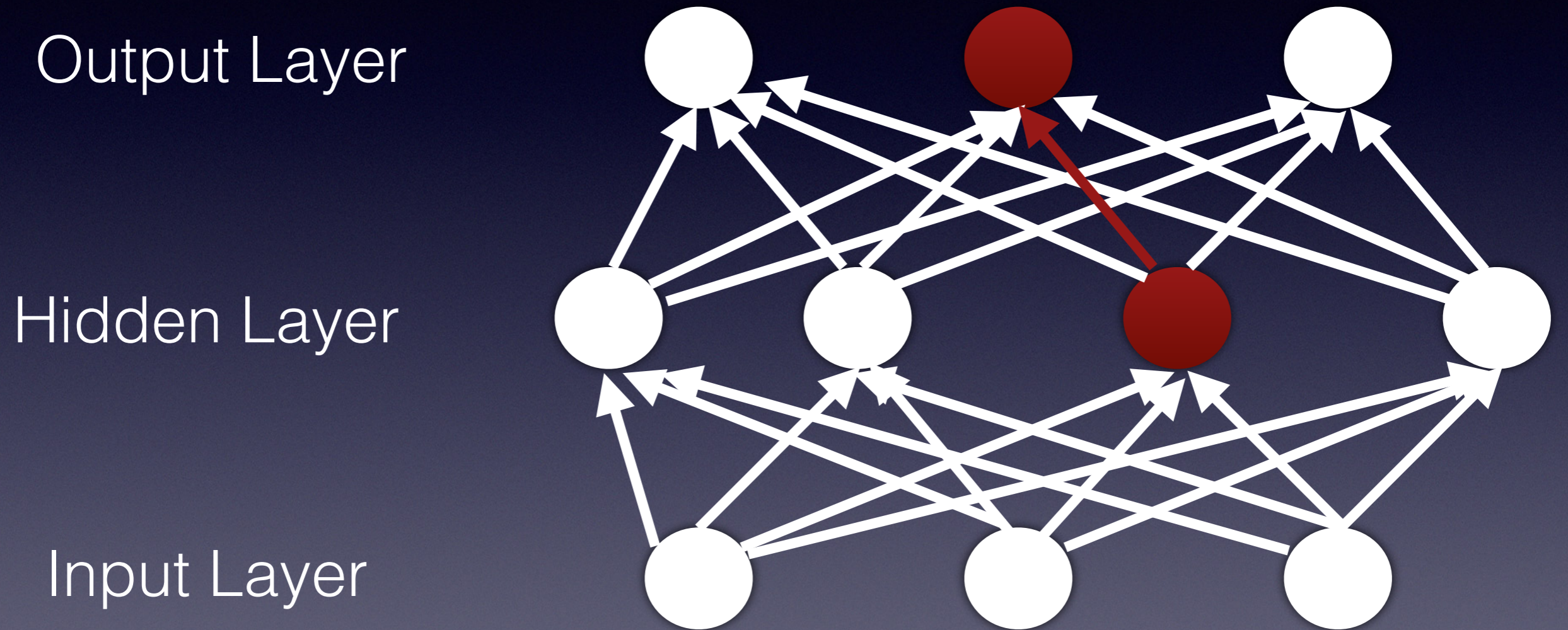
Important Neurons identified in the current layer: 2.0, 1.4, 1.5

# From Neuron to Path

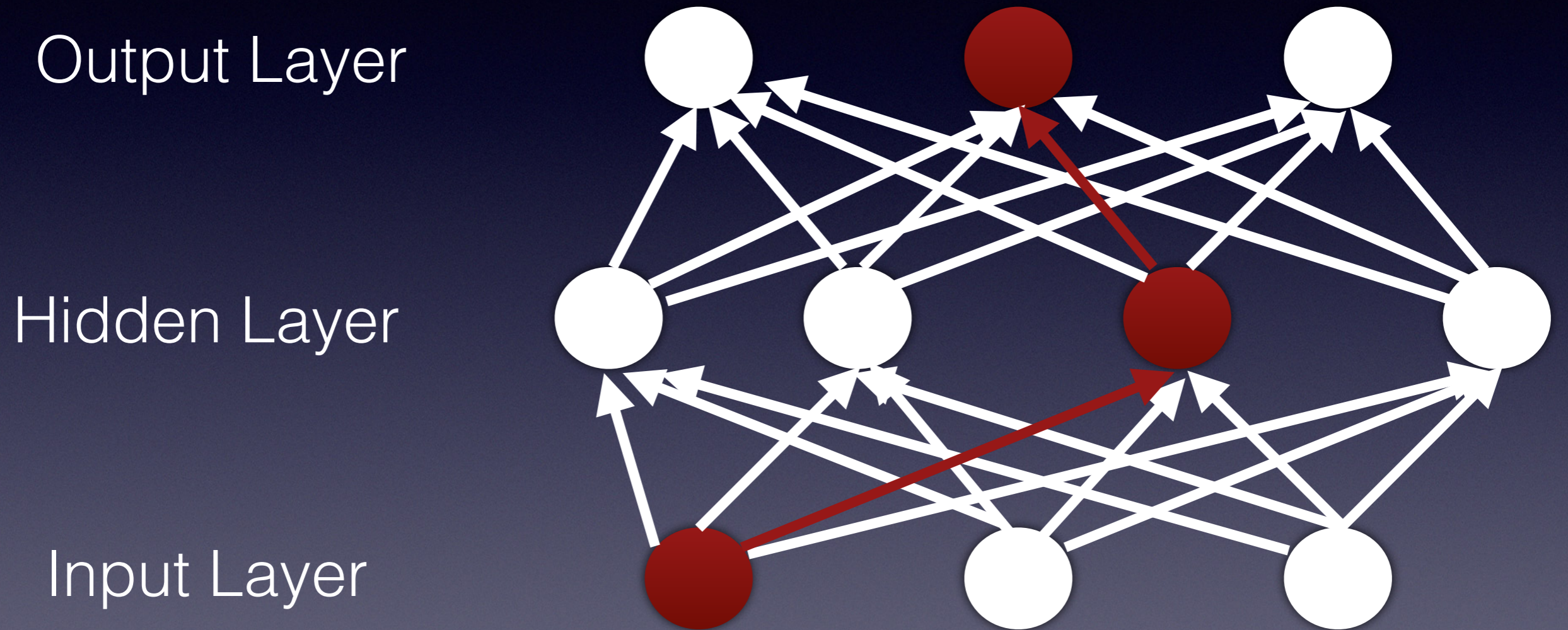




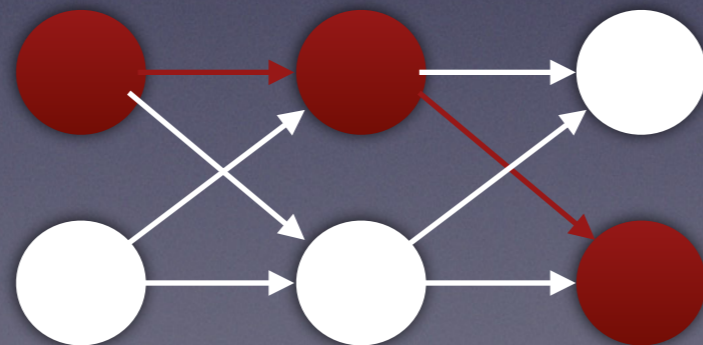
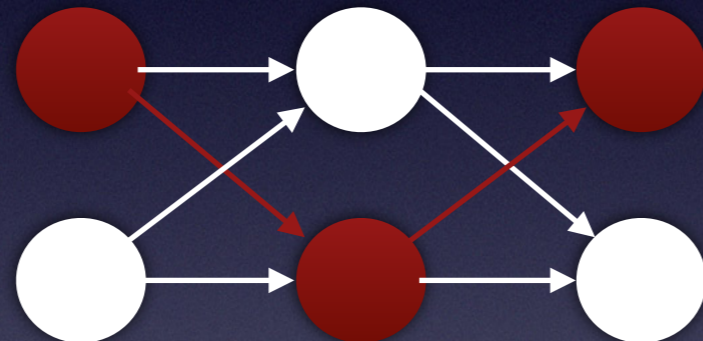
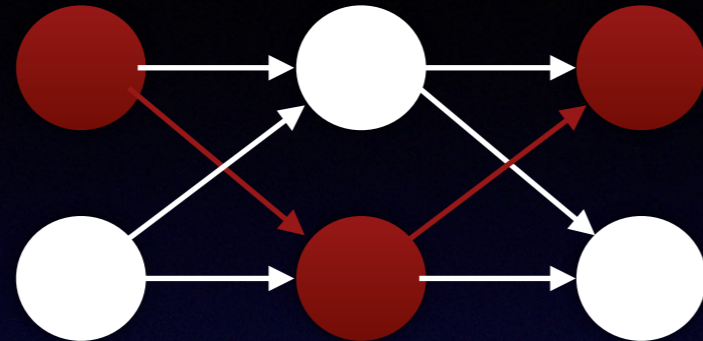
# From Neuron to Path



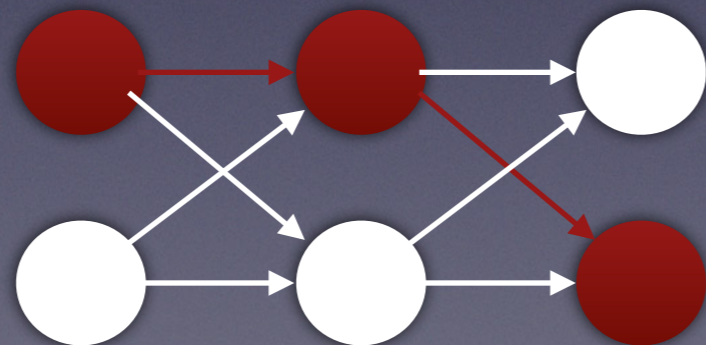
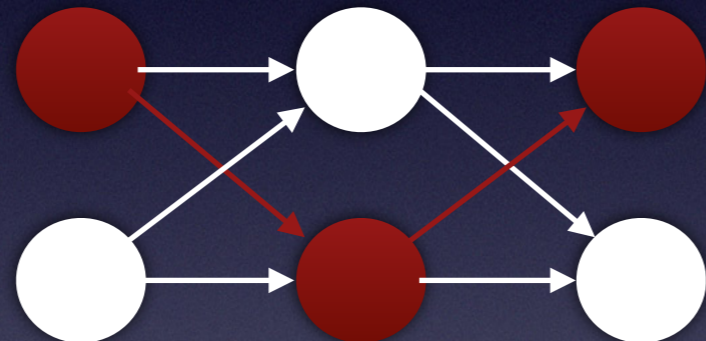
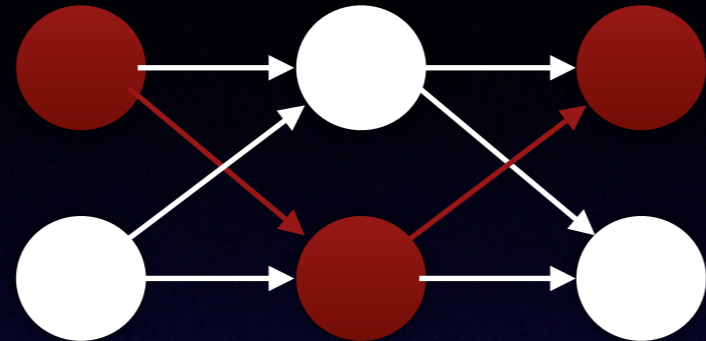
# From Neuron to Path



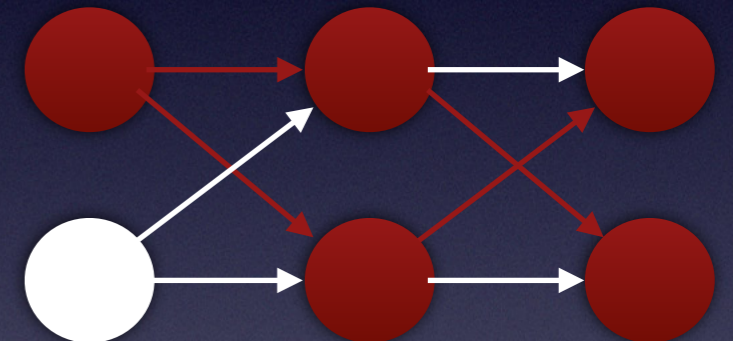
# Class Path



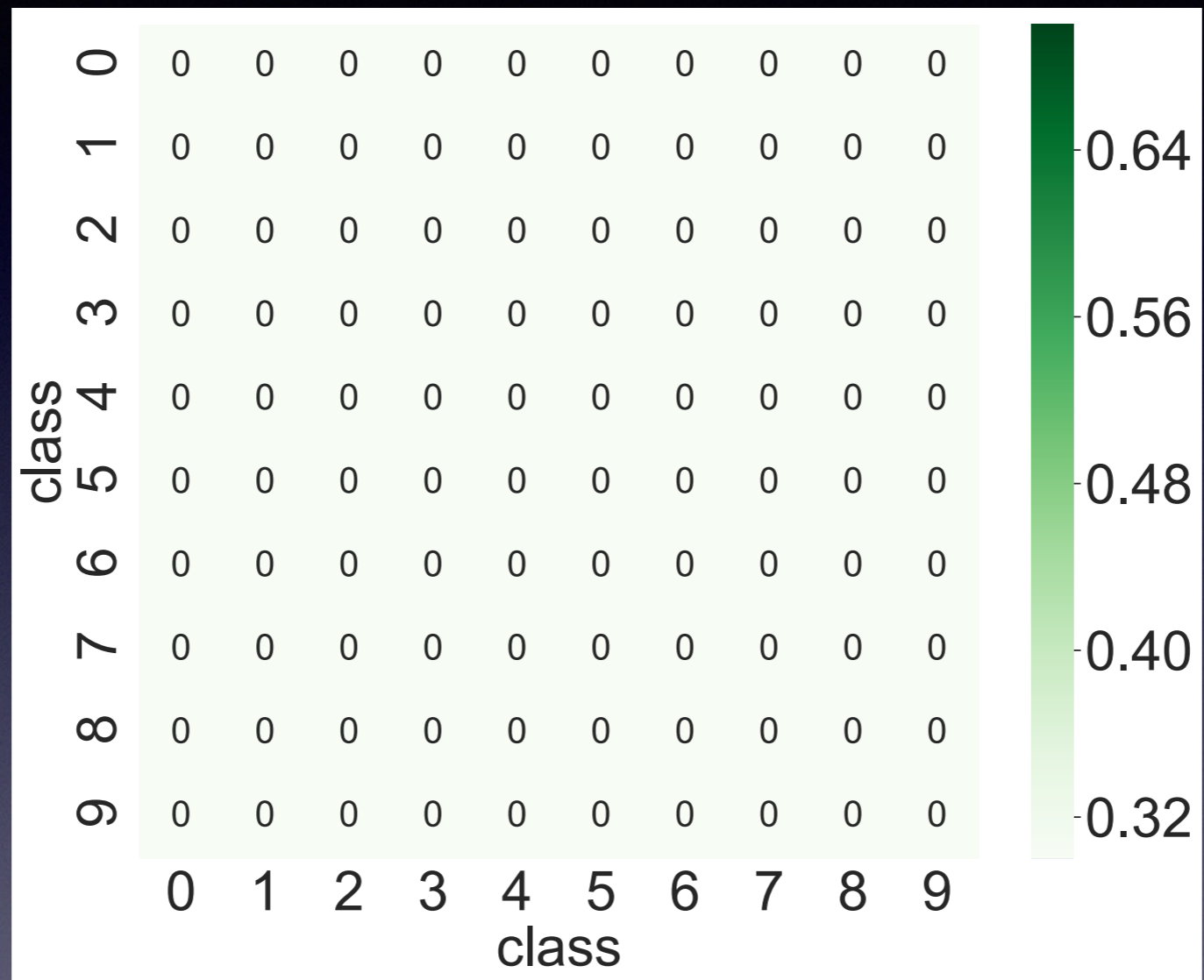
# Class Path



Union

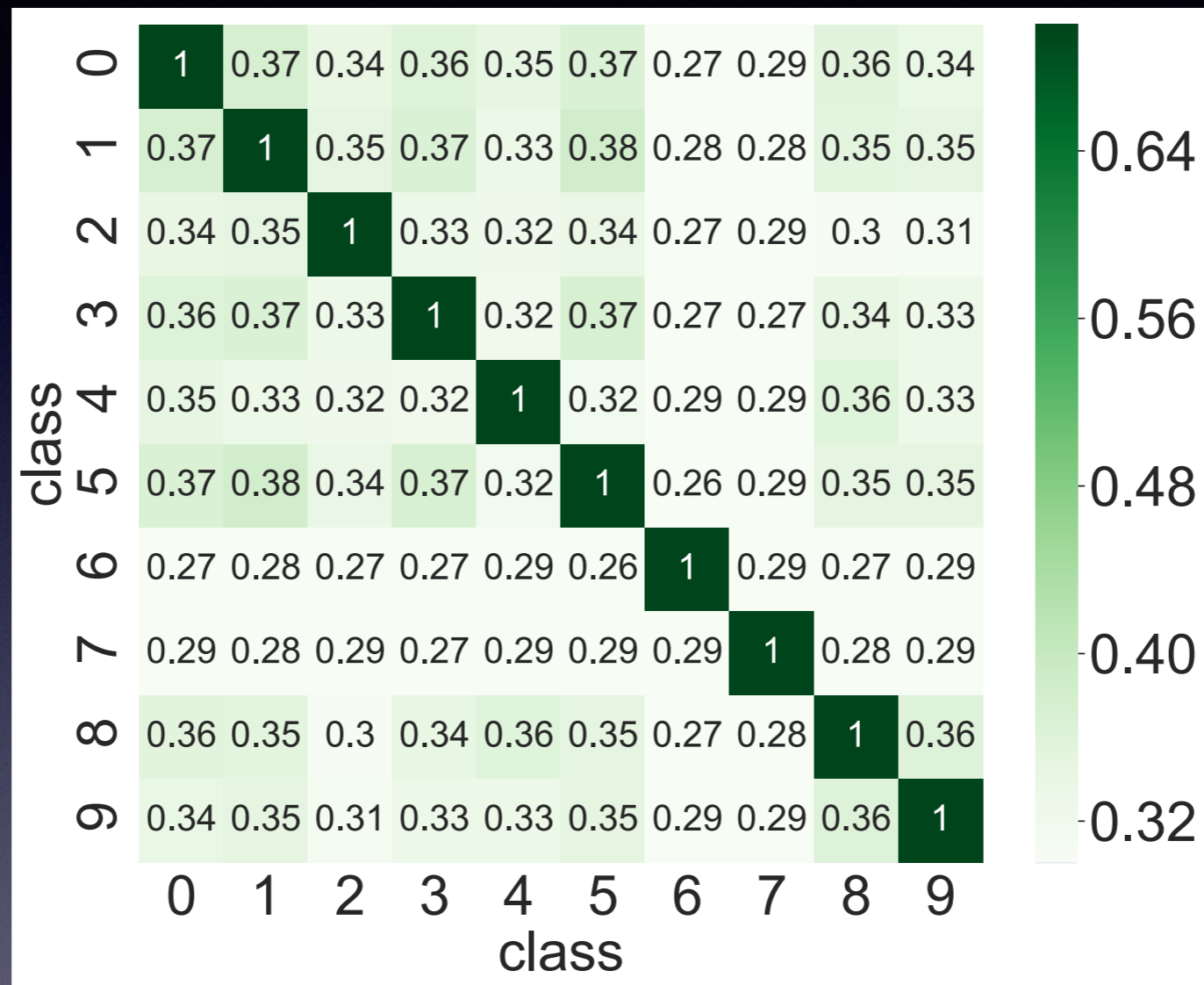


# Class Path Similarity



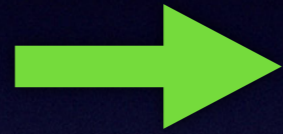
AlexNet @ ImageNet

# Class Path Similarity

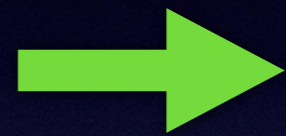


AlexNet @ ImageNet

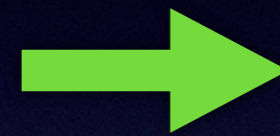
# Ptolemy Overview



# Ptolemy Overview



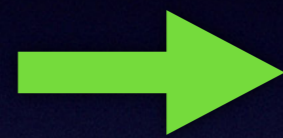
Neural  
Networks



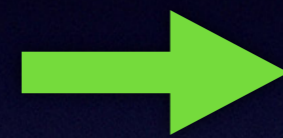
**“Cat”**



# Ptolemy Overview



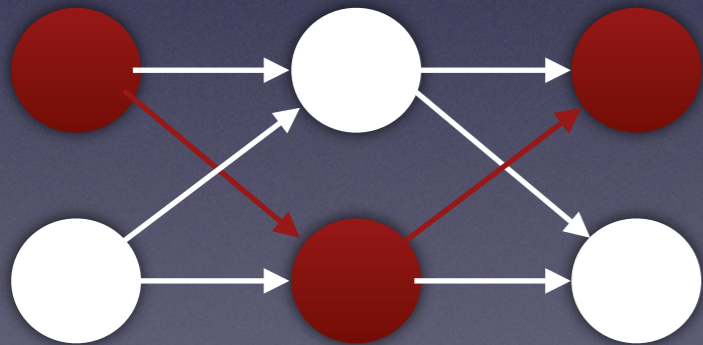
Neural  
Networks



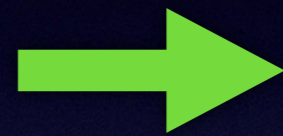
“Cat”



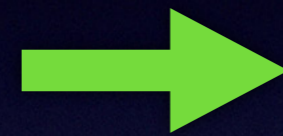
Extract



# Ptolemy Overview



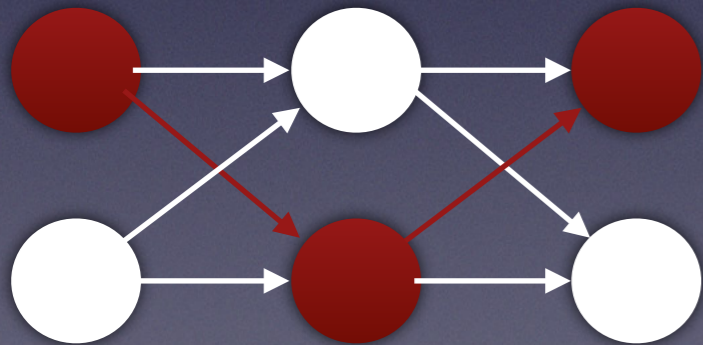
Neural  
Networks



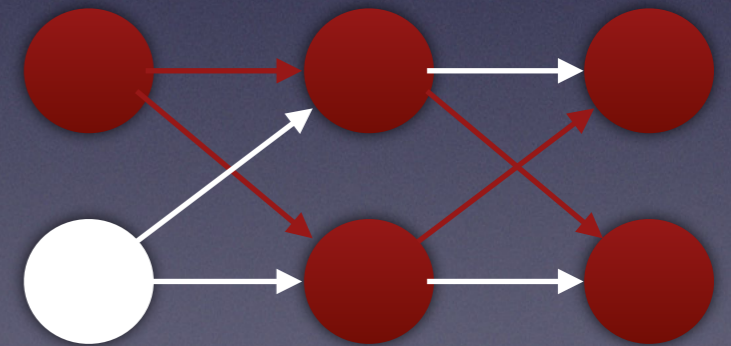
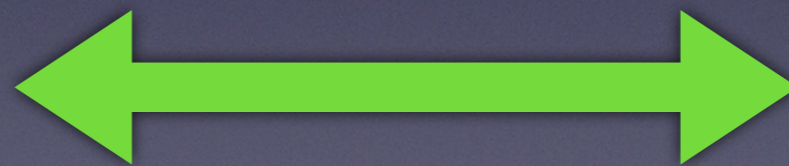
“Cat”



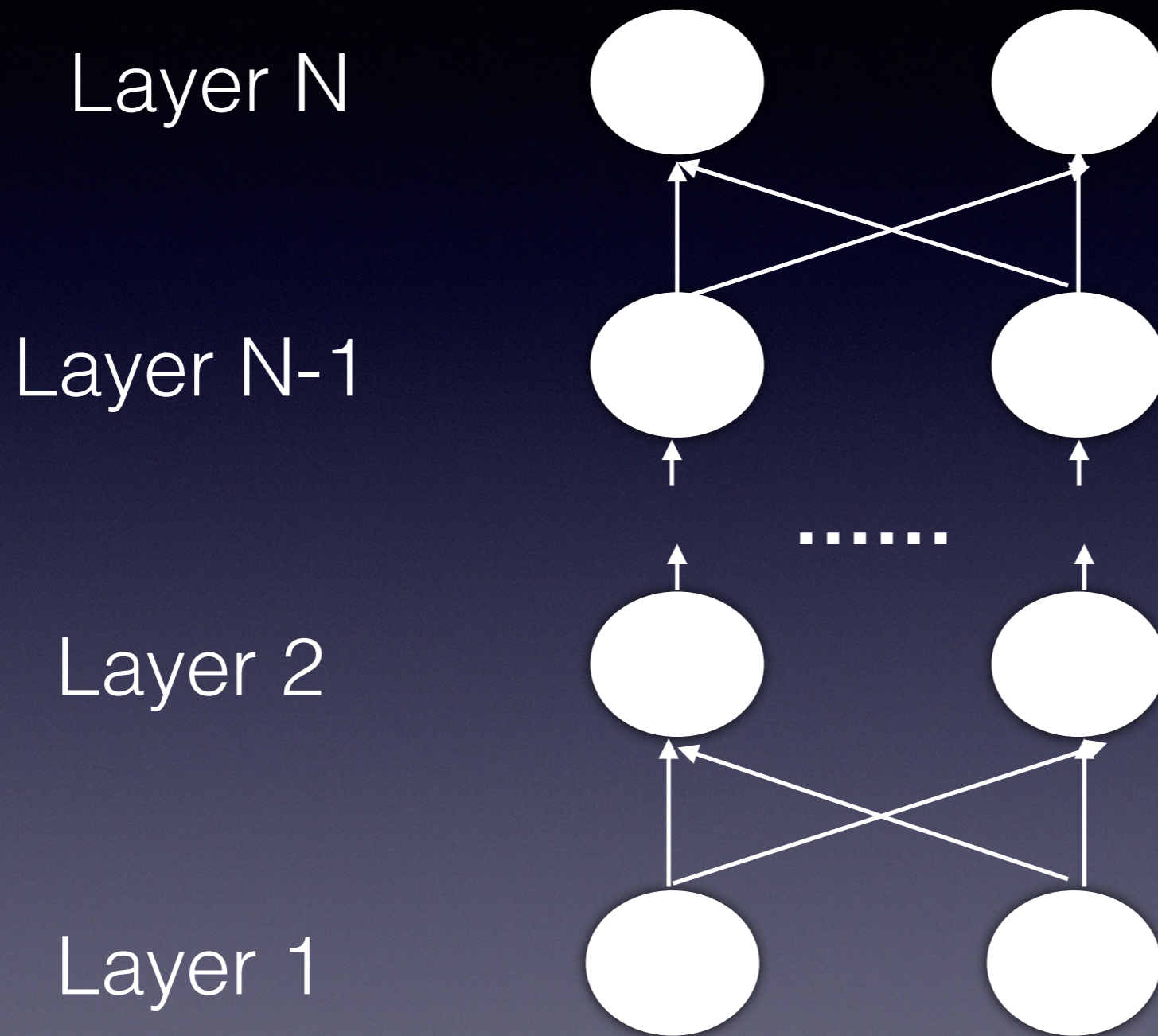
Extract



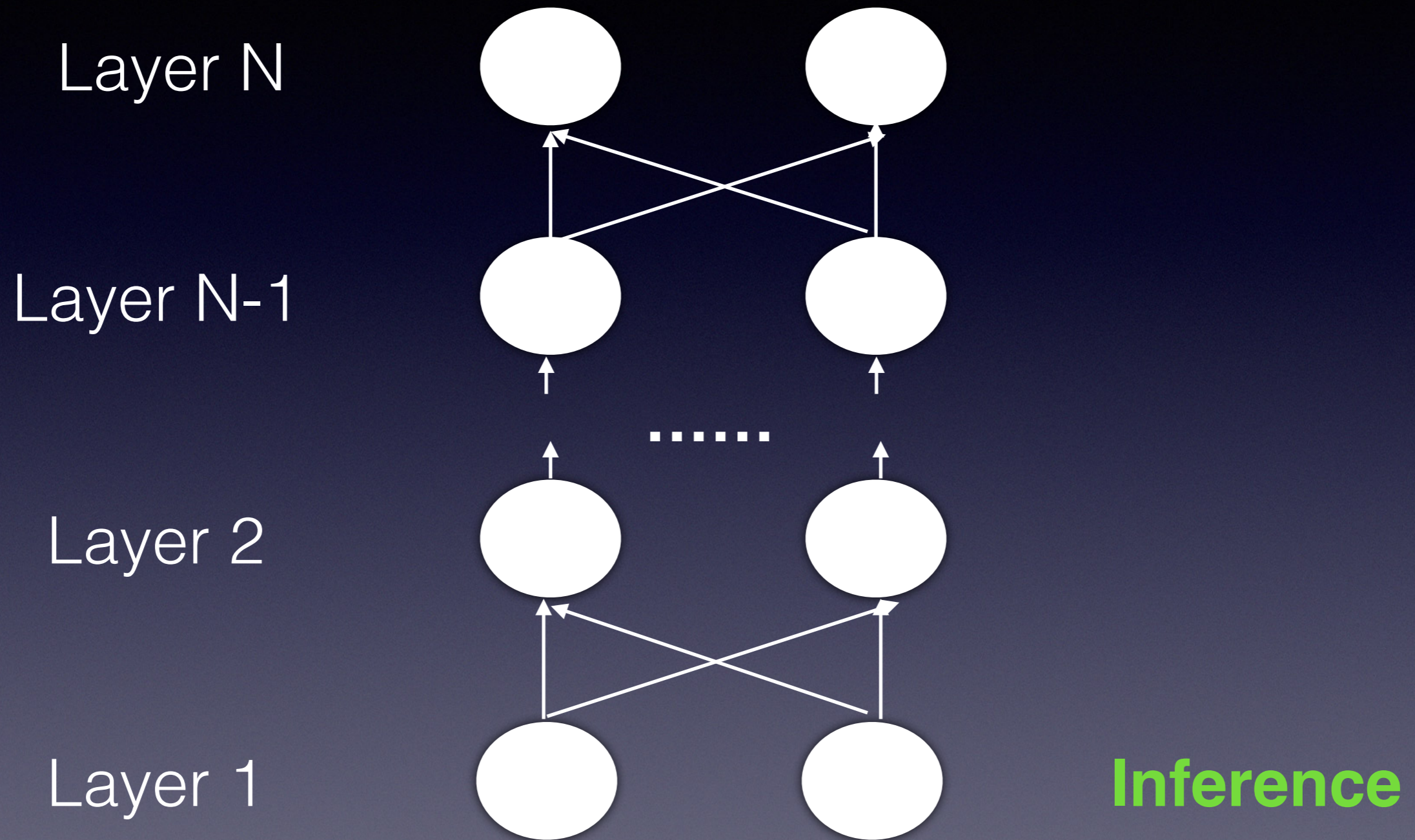
Compare



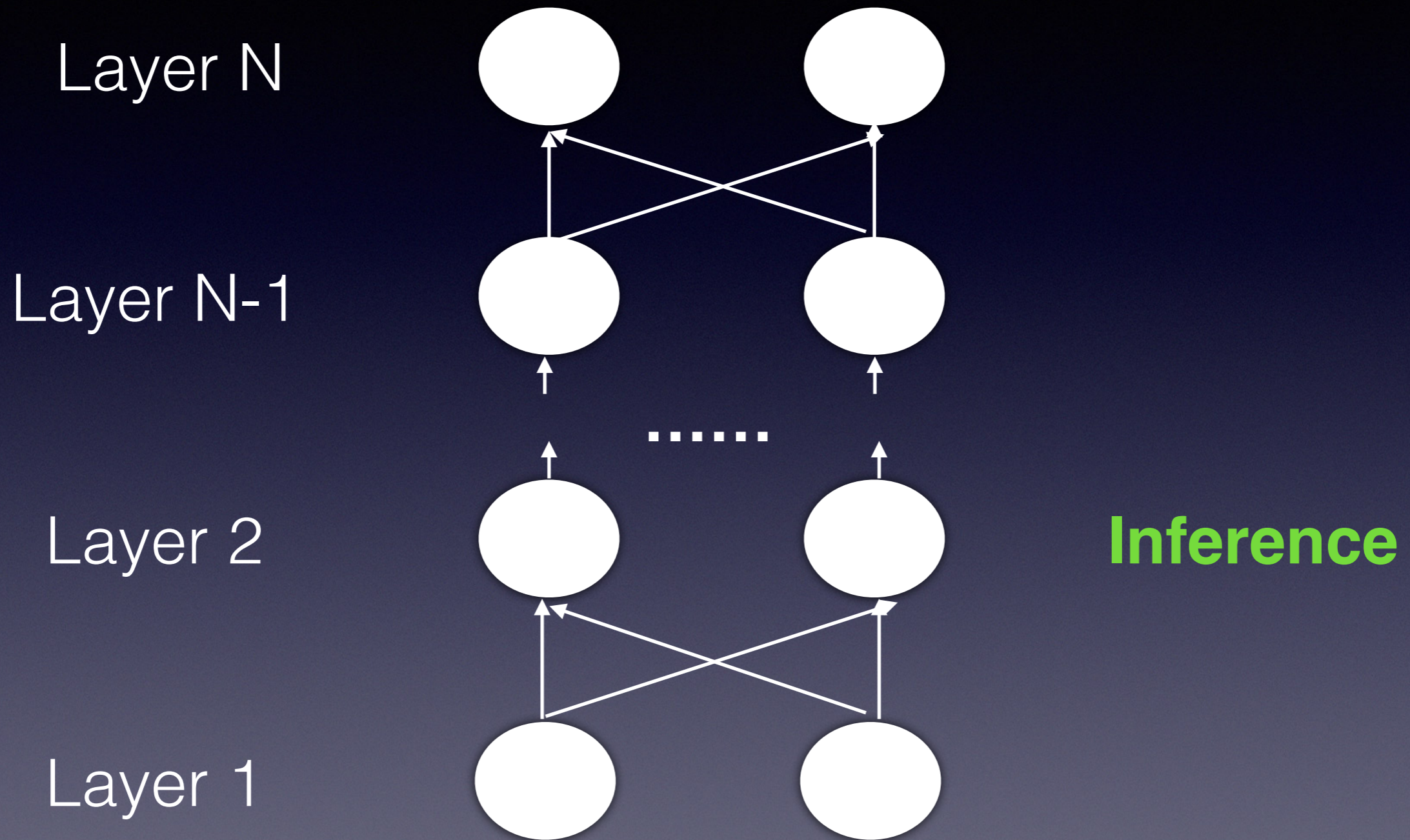
# Ptolemy Pipeline



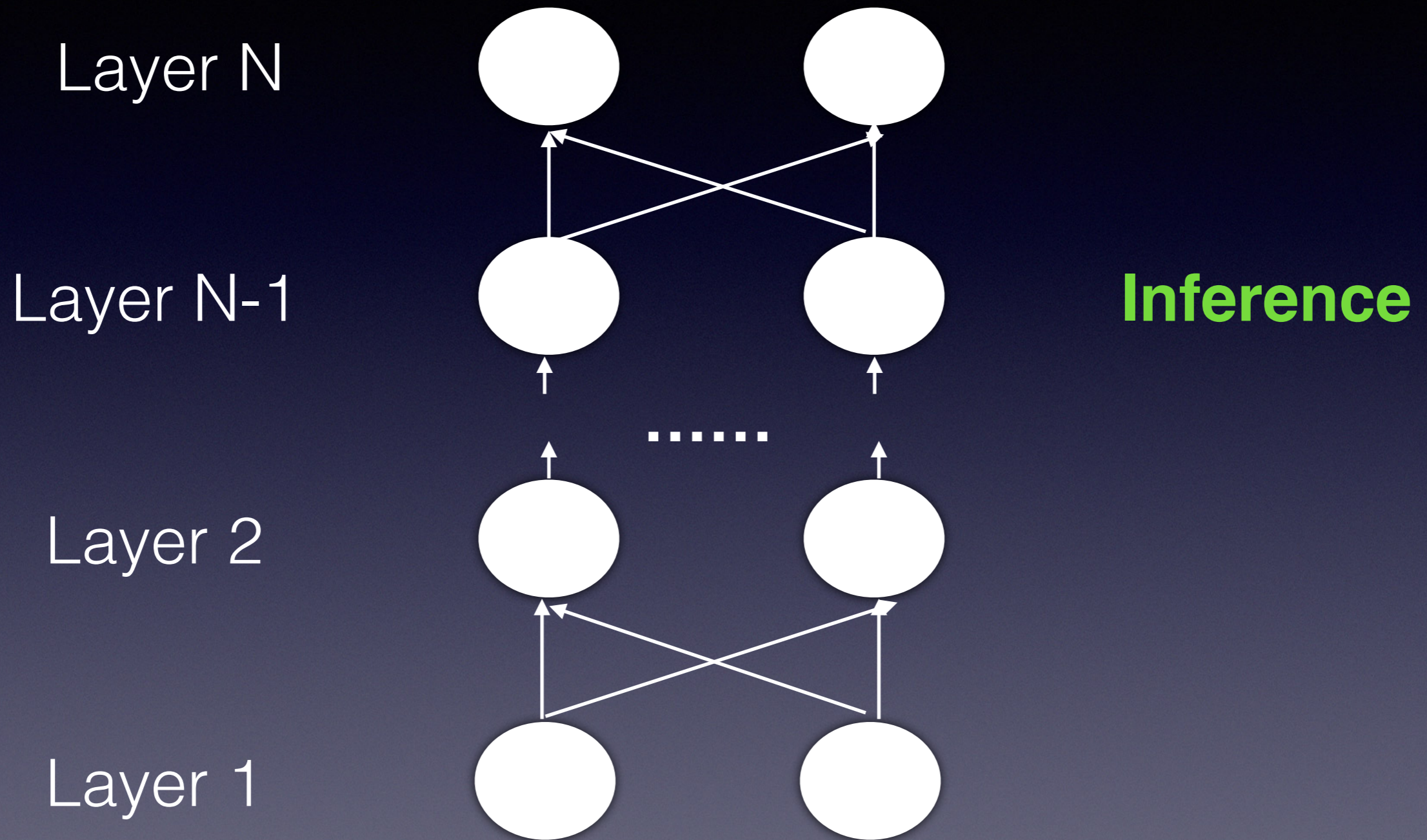
# Ptolemy Pipeline



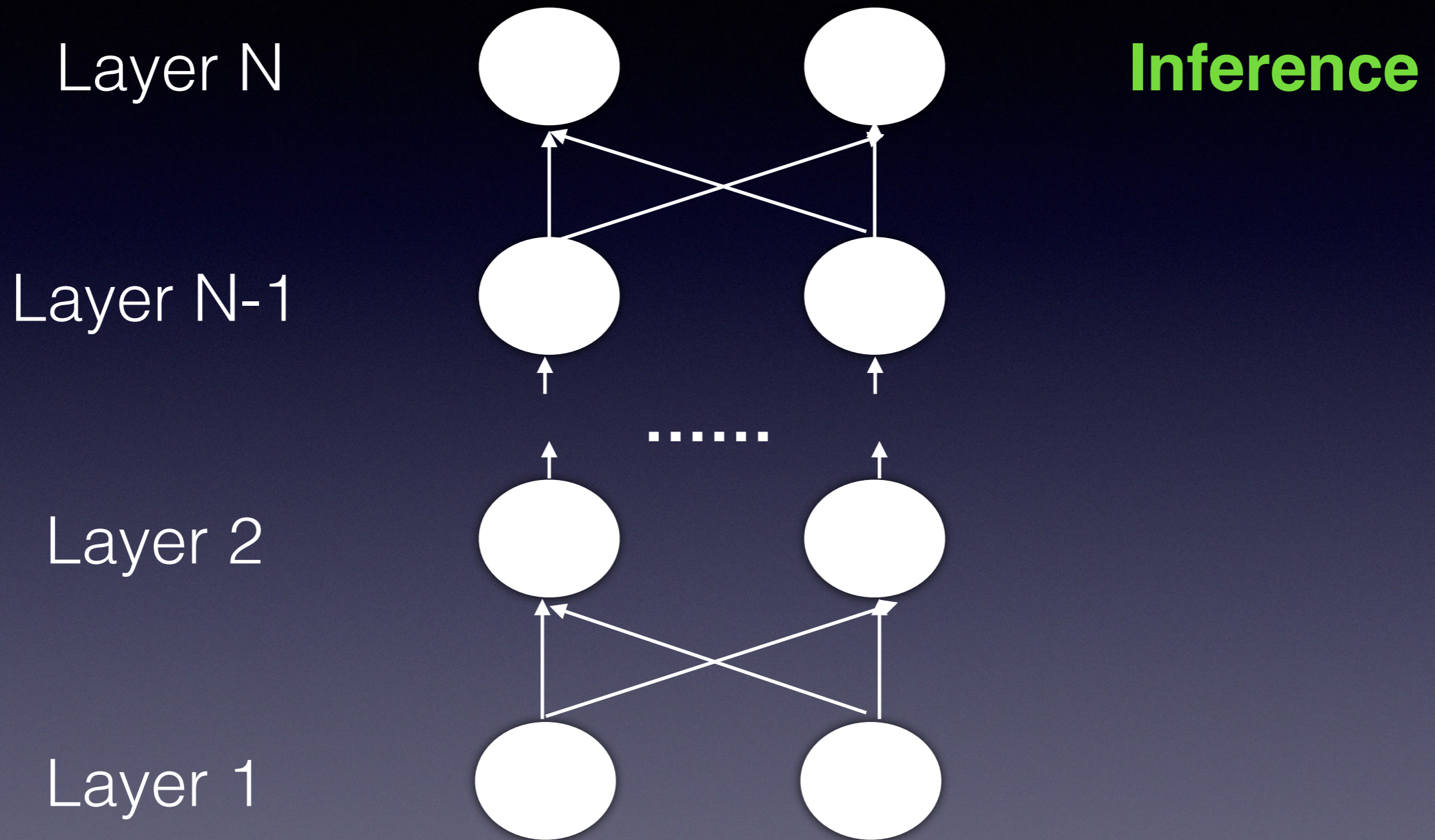
# Ptolemy Pipeline



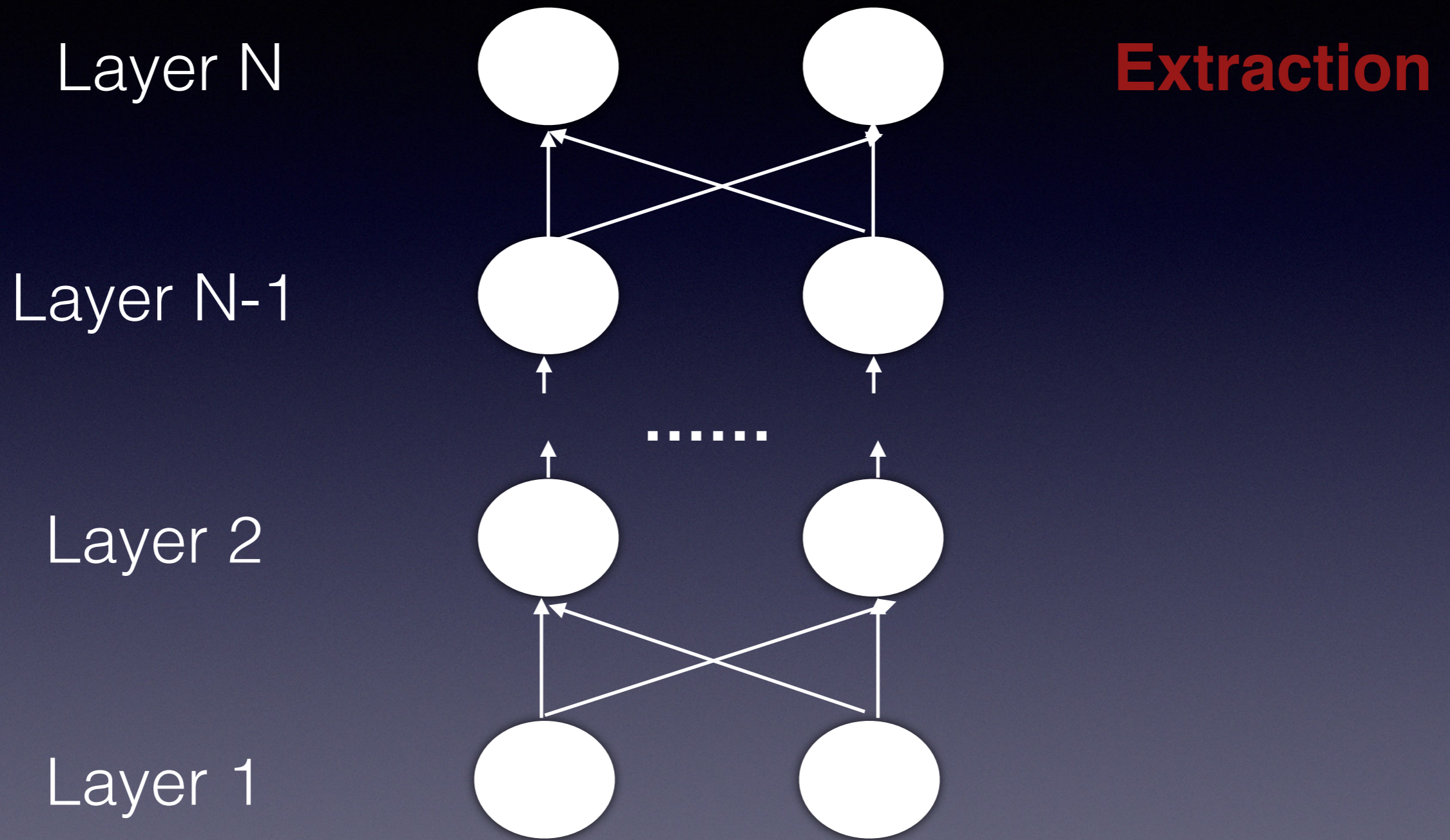
# Ptolemy Pipeline



# Ptolemy Pipeline

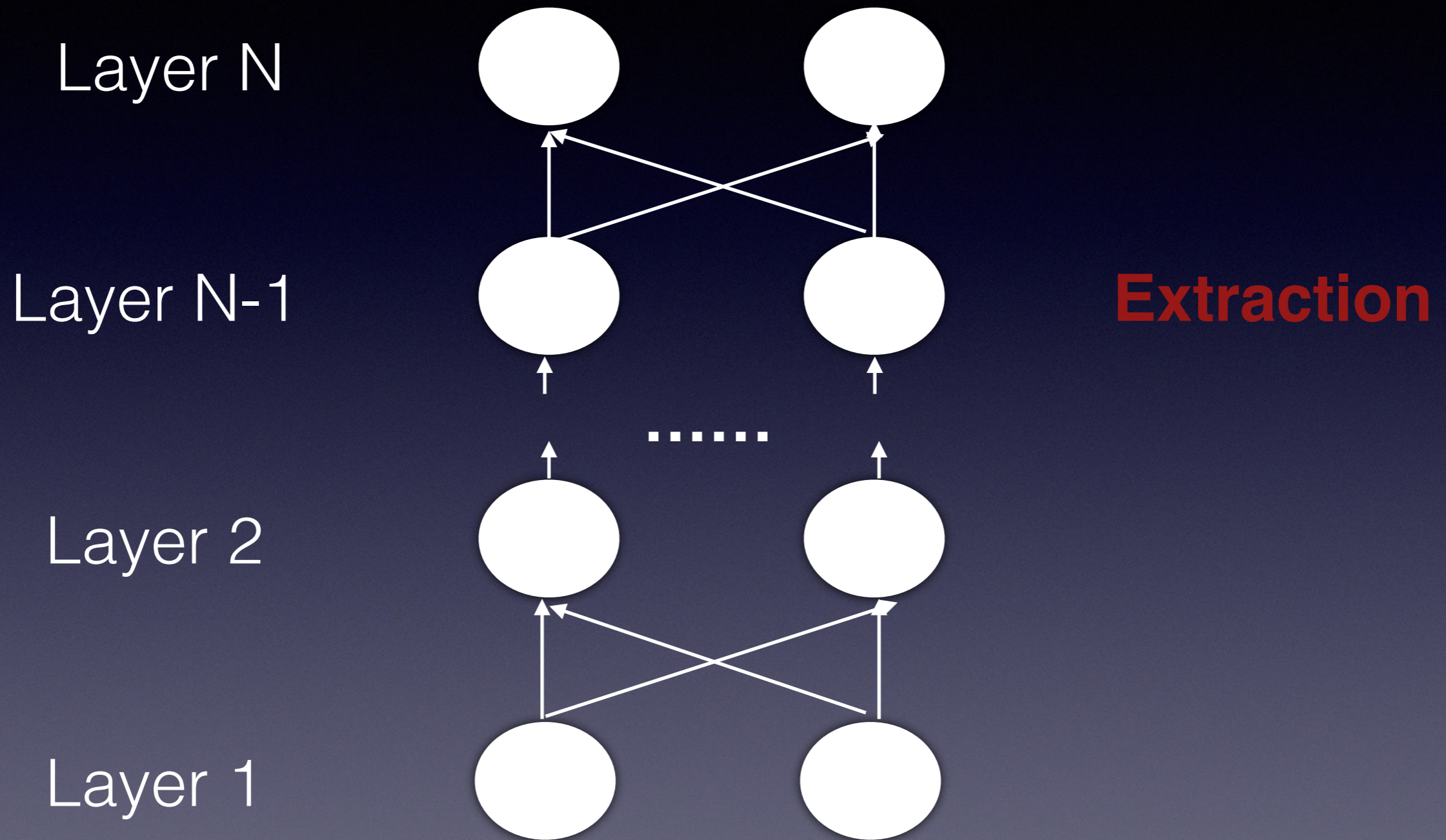


# Ptolemy Pipeline

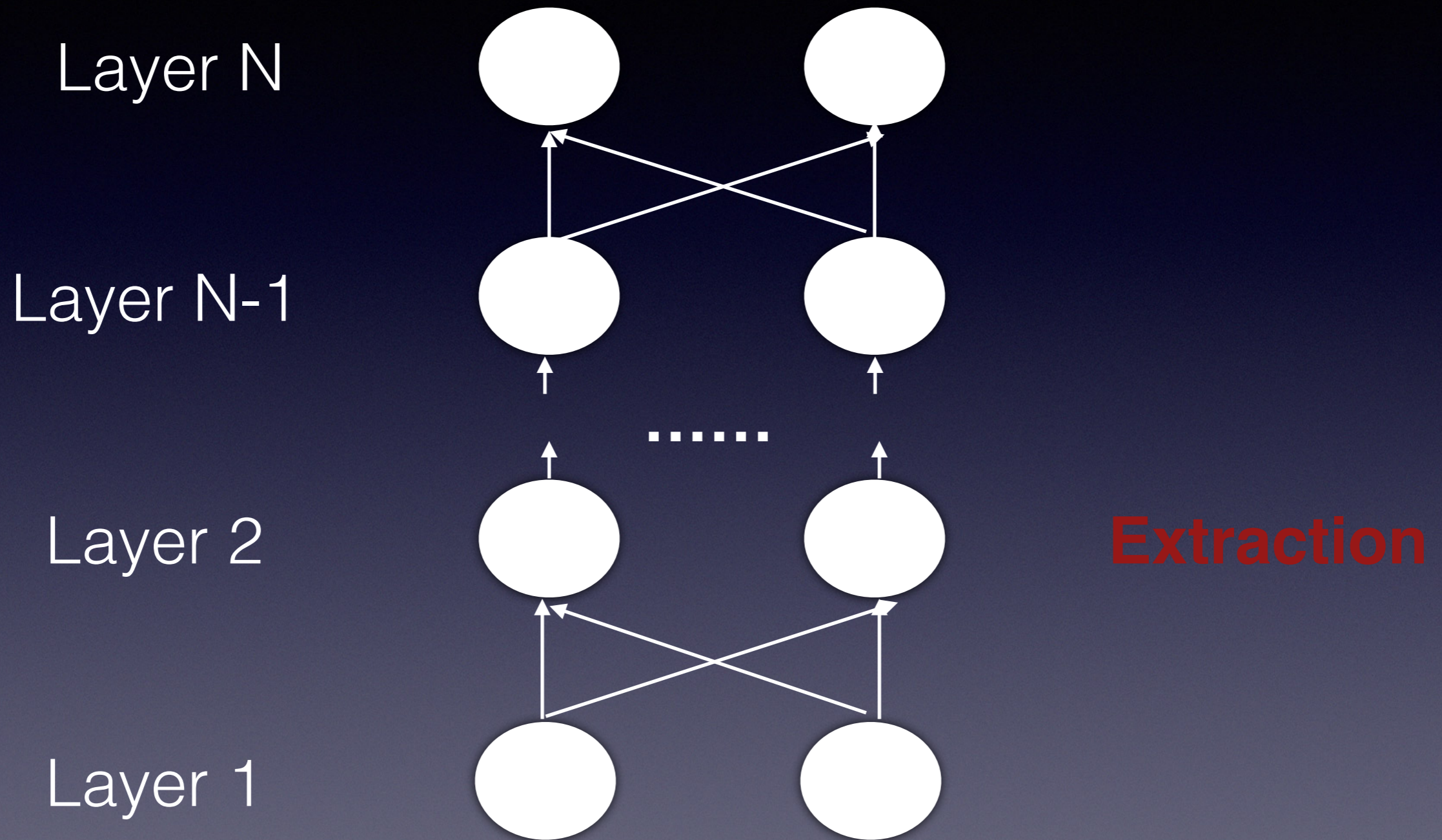




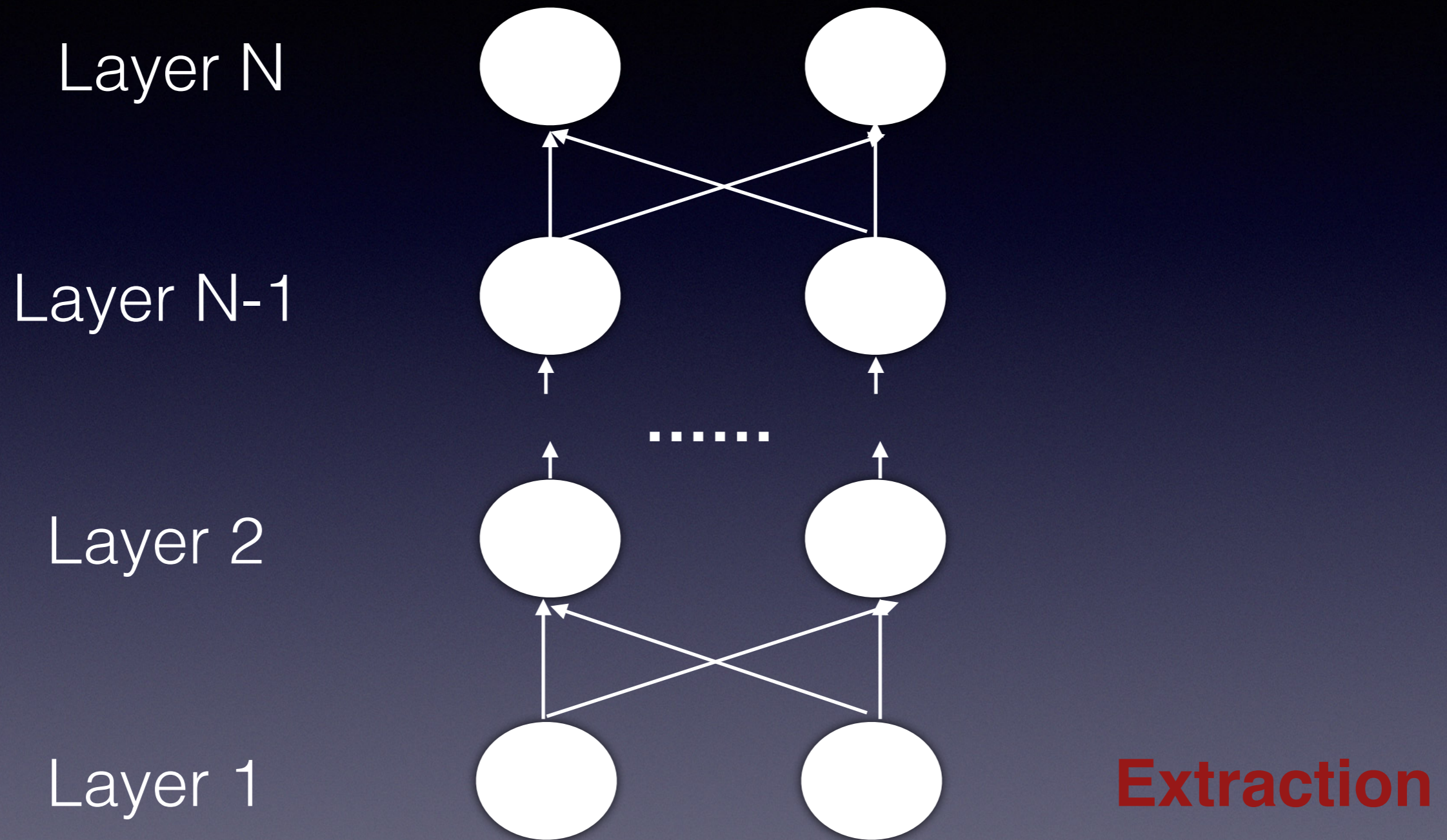
# Ptolemy Pipeline



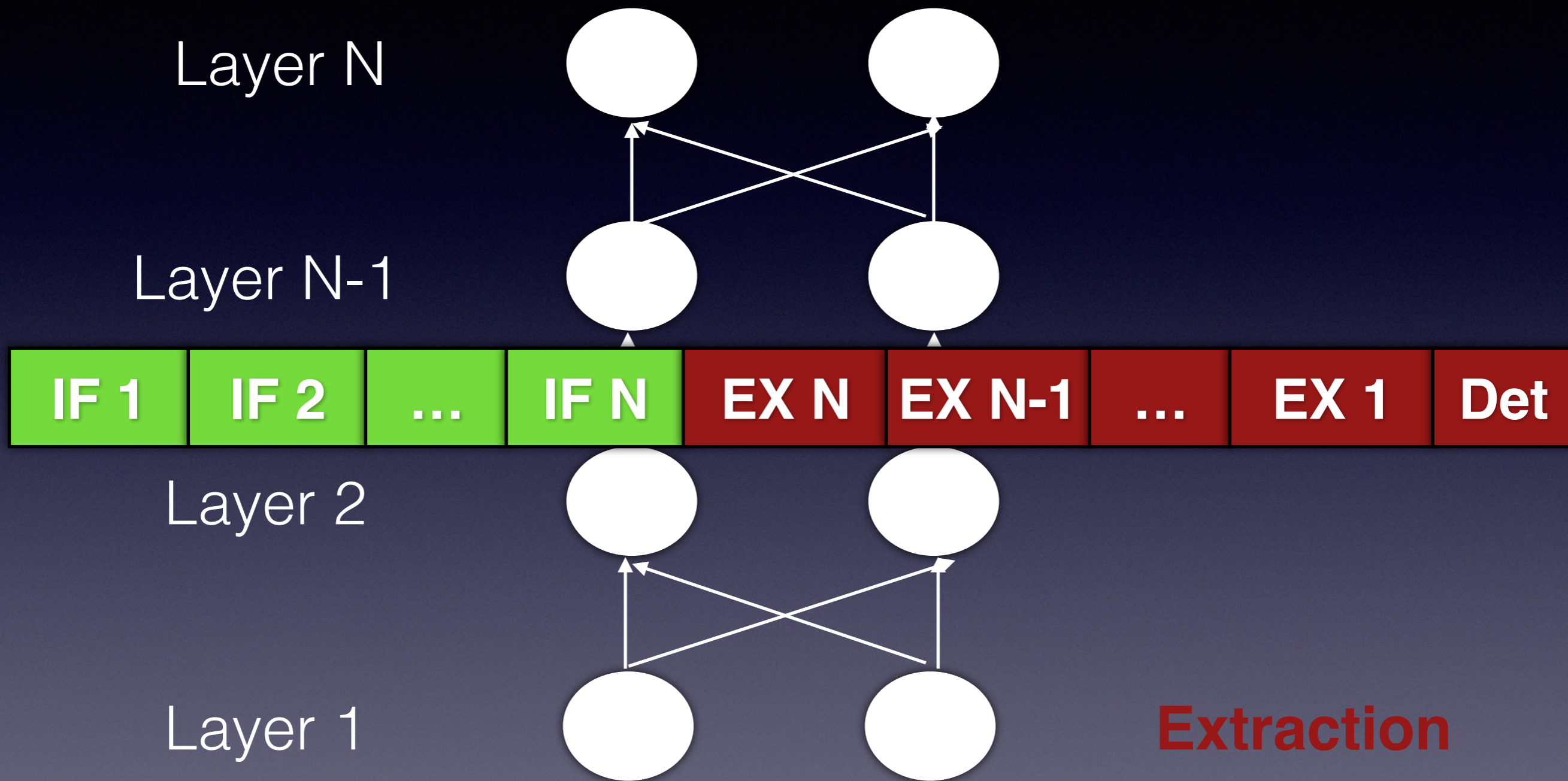
# Ptolemy Pipeline



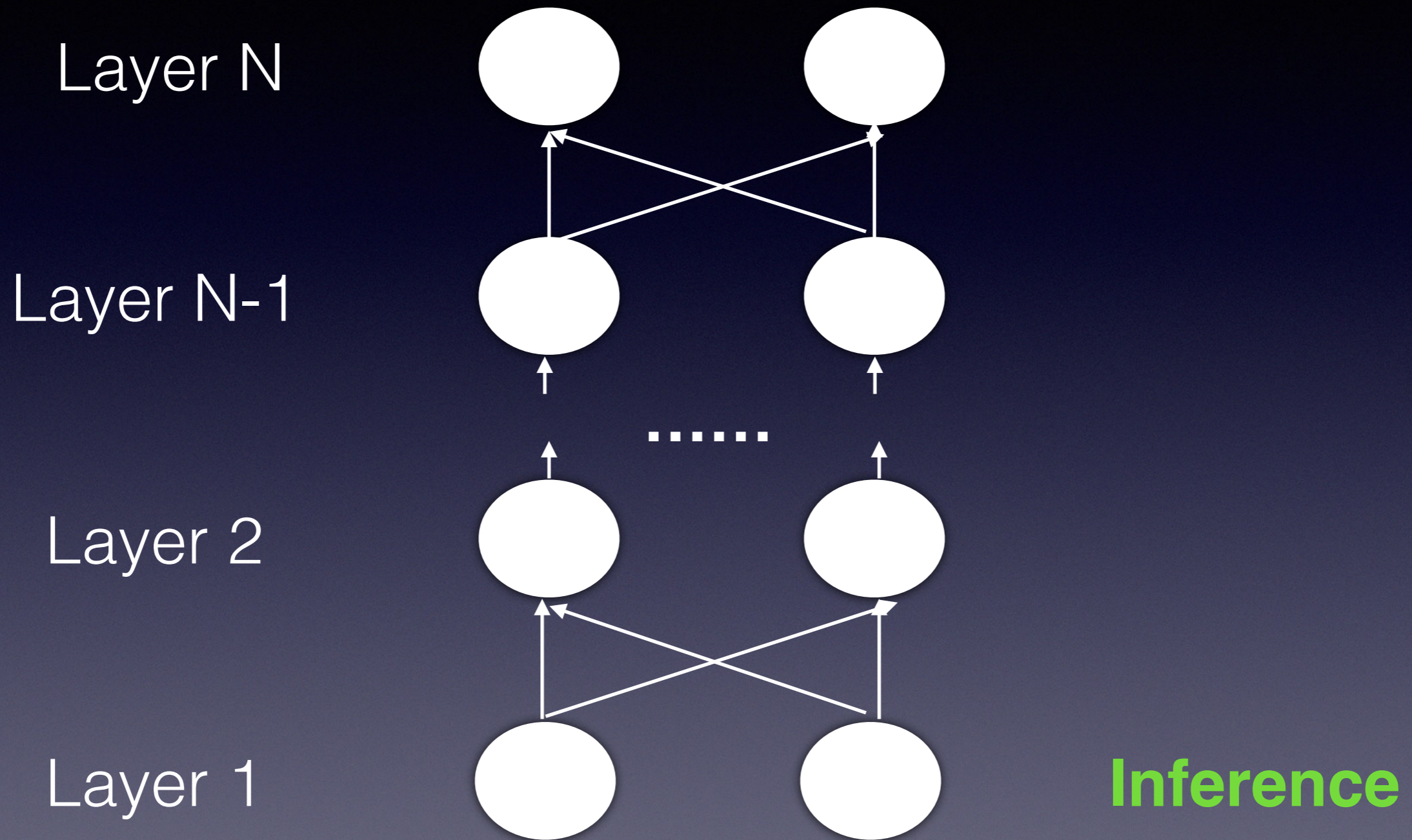
# Ptolemy Pipeline



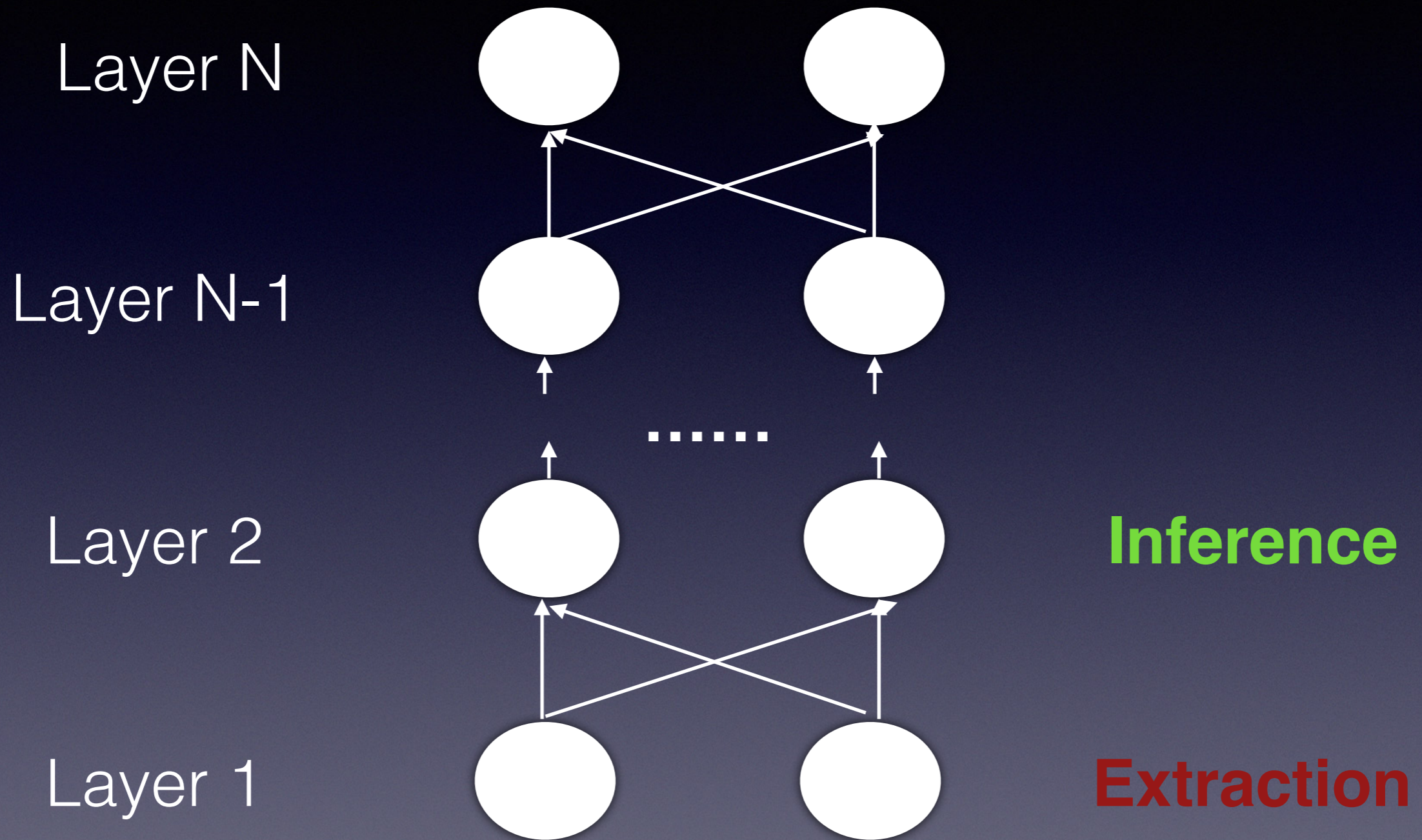
# Ptolemy Pipeline



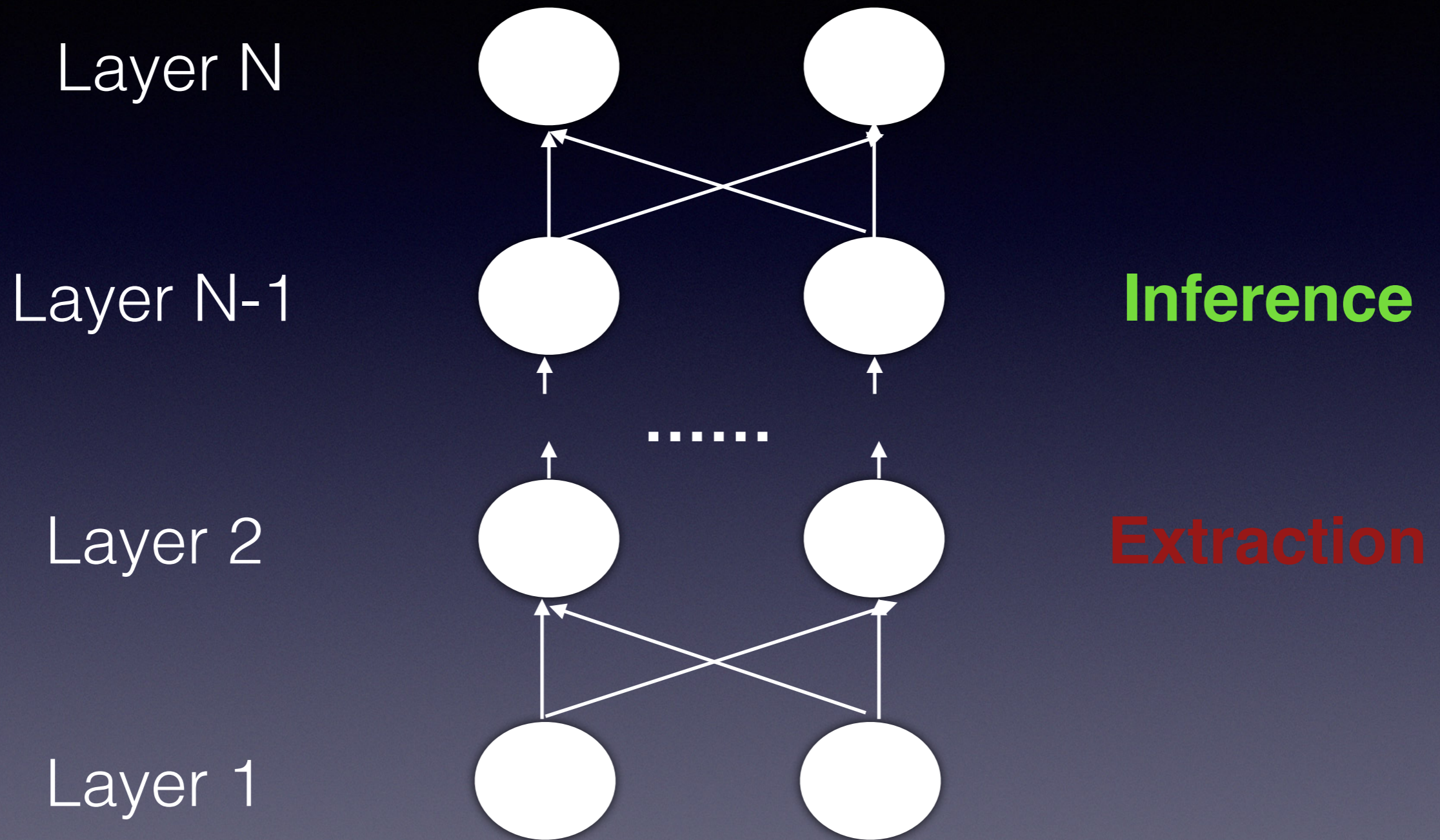
# Algorithmic Variation



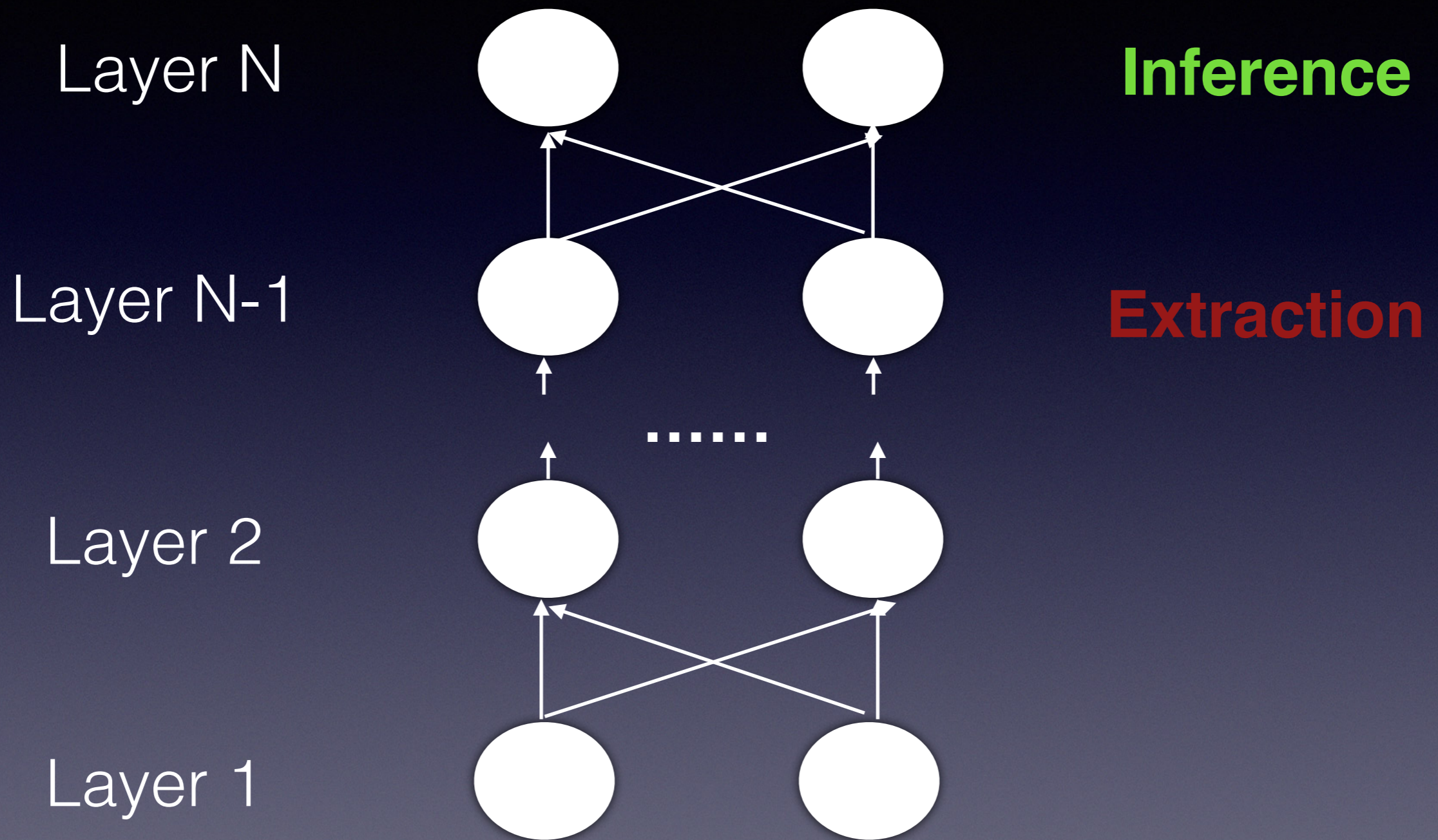
# Algorithmic Variation



# Algorithmic Variation

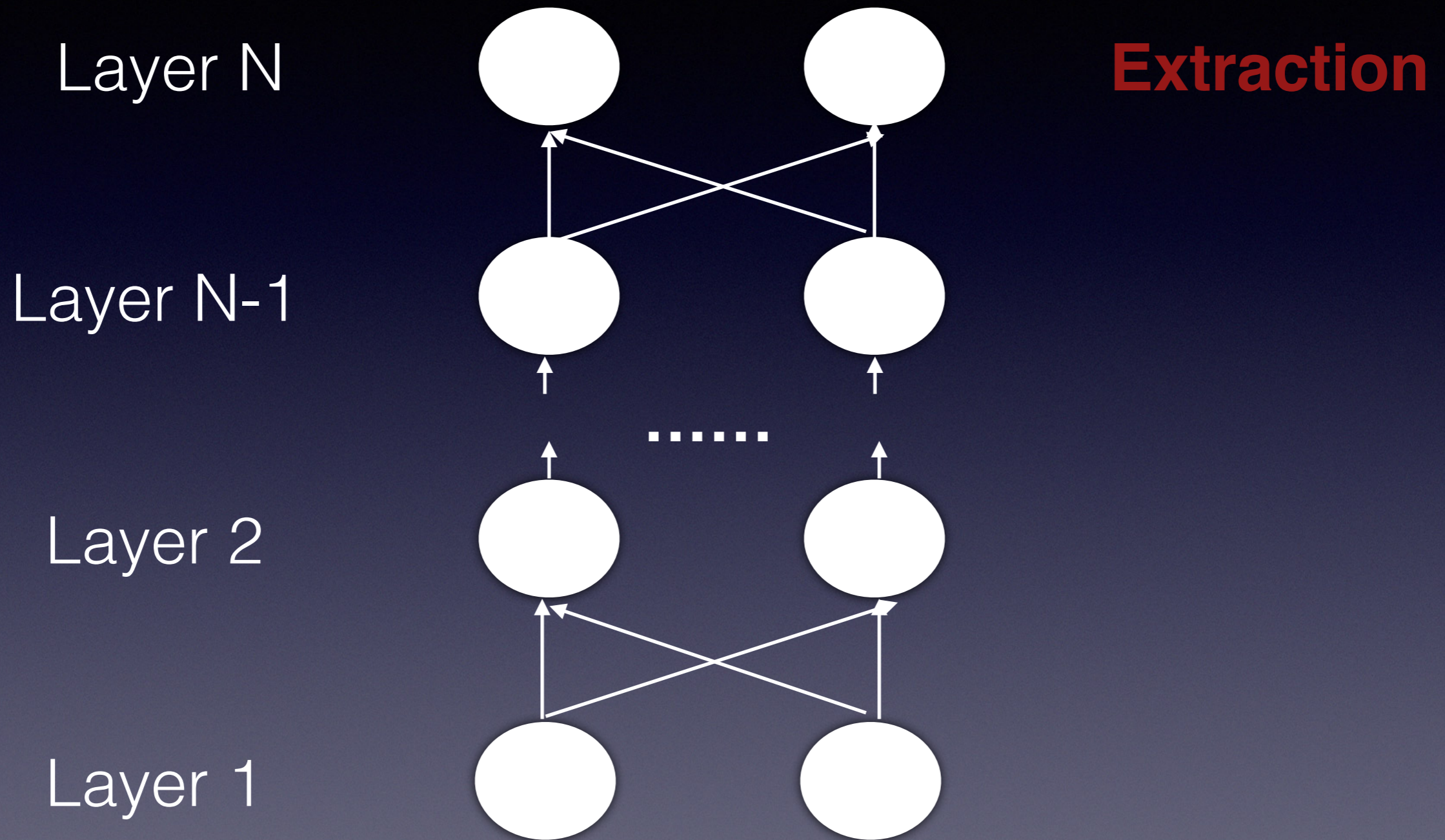


# Algorithmic Variation

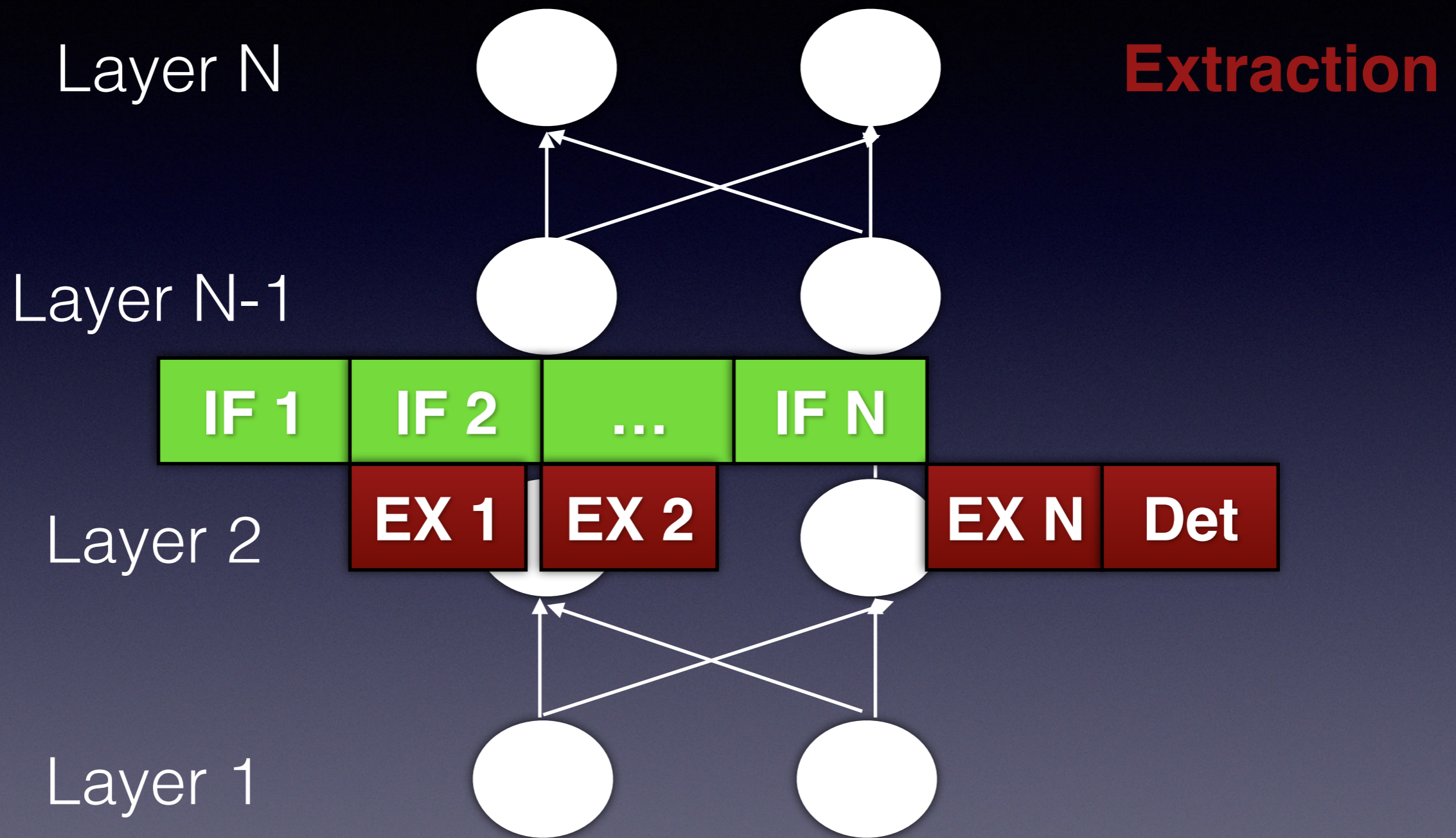




# Algorithmic Variation



# Algorithmic Variation



# Algorithmic Variation

Sorting



Threshold



IF: Inference, EX: Extraction, Det: Detection

# Algorithmic Variation

Full Extraction



Partially Extraction



IF: Inference, EX: Extraction, Det: Detection

# Framework

**Backward**

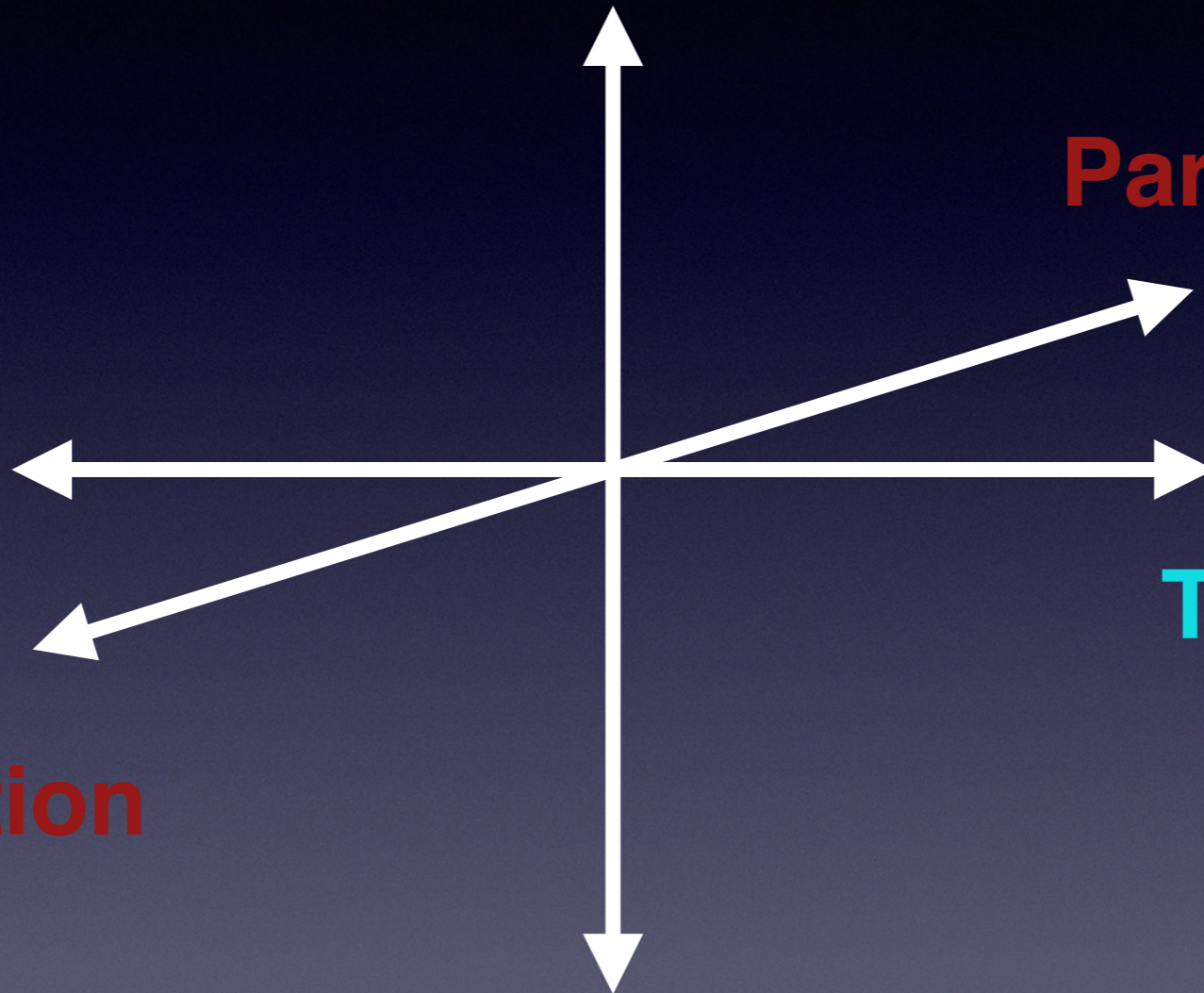
**Partially Extraction**

**Sorting**

**Thresholding**

**Fully Extraction**

**Forward**



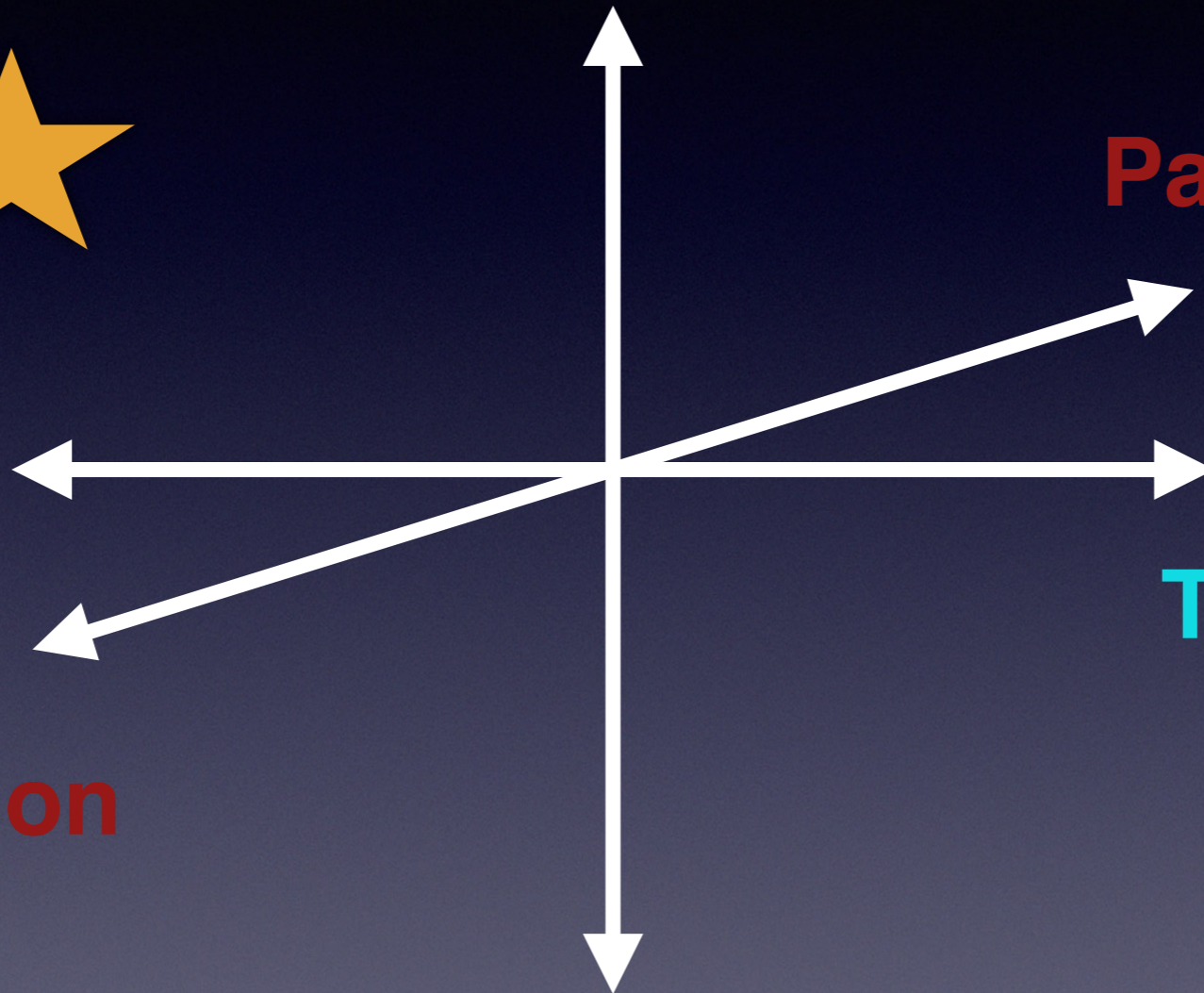
# Framework

**Backward**



**Partial Extraction**

**Sorting**



**Thresholding**

**Full Extraction**

**Forward**



**= Backward + Fully Extraction + Sorting**

# Interface

- High-level: Python-based, user define input

# Interface

- High-level: Python-based, user define input



Compiler



- Low-level: Customized ISA



# Compiler Optimization: Layer Level

```
for j = 1 to L {  
    inf(j)  
    <extraction on layer j>  
}
```

# Compiler Optimization: Layer Level

```
for j = 1 to L {  
    inf(j)  
    <extraction on layer j>  
}
```



```
inf(1)  
for j = 1 to L {  
    inf (j+1)  
    <extraction on layer j>  
}  
<extraction on layer L>
```

# Compiler Optimization: Neuron Level

```
for j = 1 to N {  
    sort(i)  
    acum(i)  
}
```

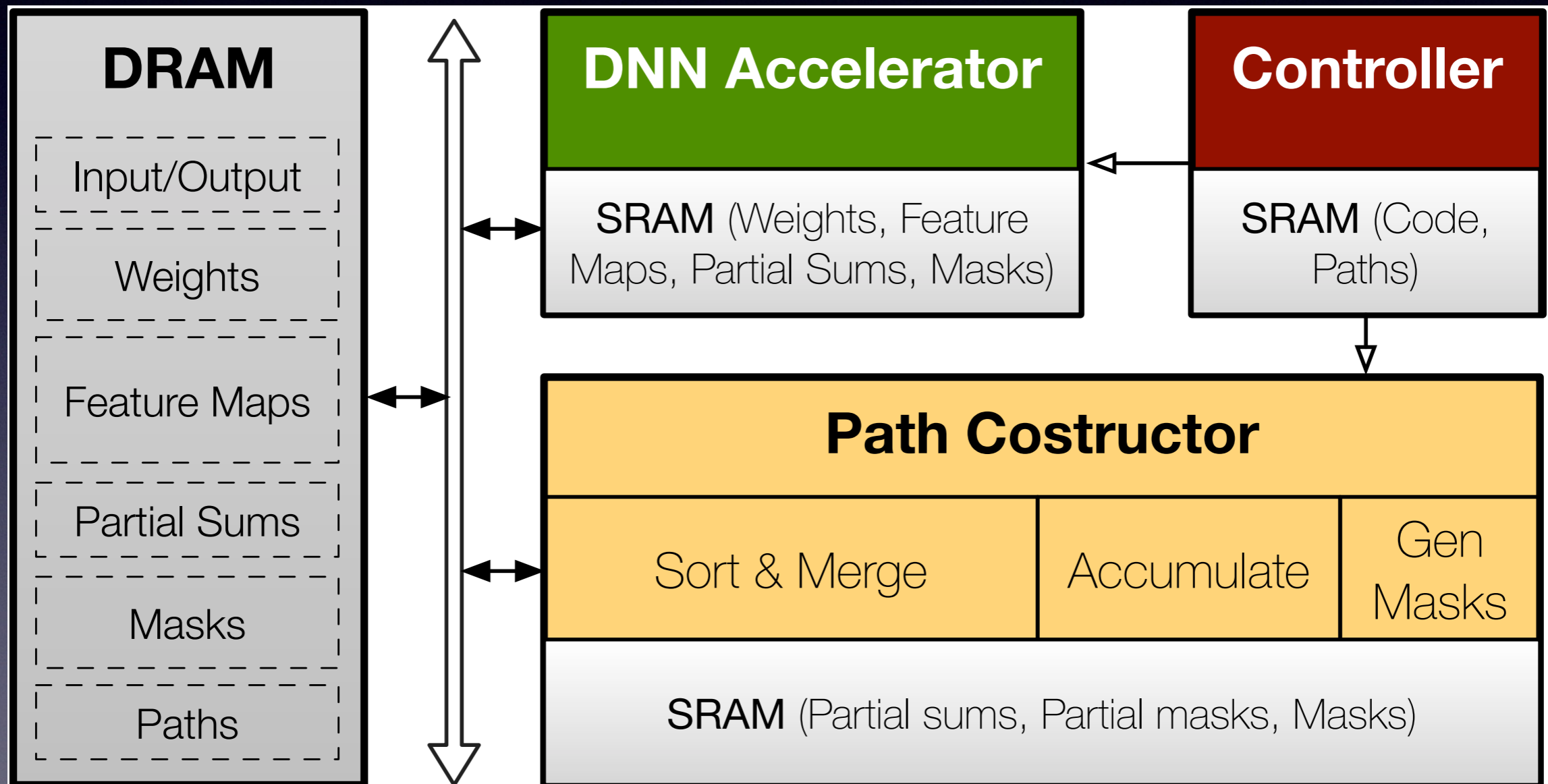
# Compiler Optimization: Neuron Level

```
for j = 1 to N {  
    sort(i)  
    acum(i)  
}
```

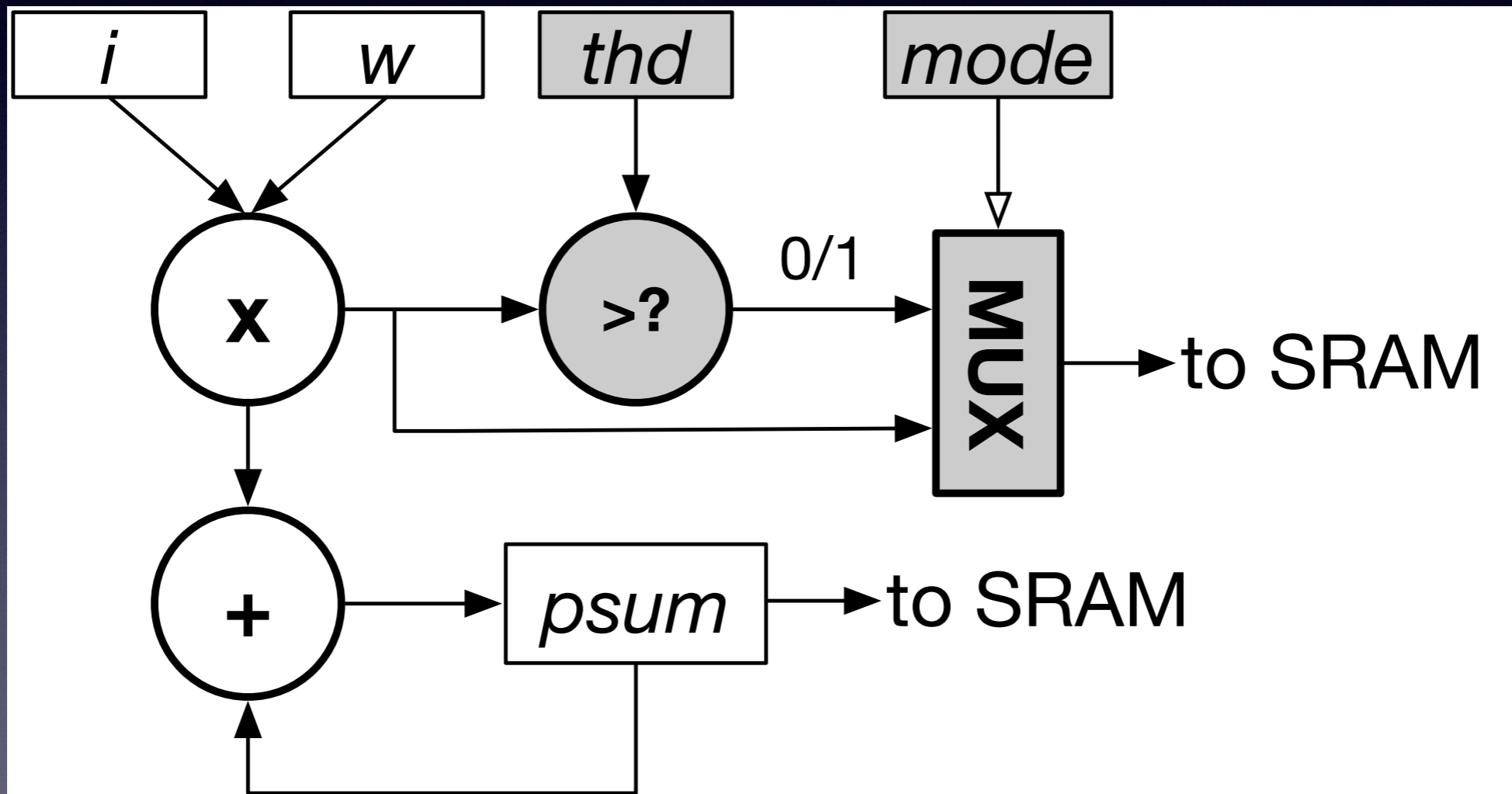


```
sort(1)  
for i = 1 to N-1 {  
    sort(i+1)  
    acum(i)  
}  
acum(N)
```

# Architecture Overview



# Enhanced MAC unit



# Evaluation

Network	AlexNet, ResNet
Dataset	Cifar10, Cifar100, ImageNet
Attacks	BIM, CWL2, DeepFool, FGSM, JSMA
Adaptive Attacks	Self constructed
Baselines	EP[1], CDRP[2]

[1]Y. Qiu, J. Leng, C. Guo, et.al, Adversarial Defense Through Network Profiling Based Path Extraction

[2]Y. Wang, H. Su, B. Zhang, X. Hu, Interpret neural networks by identifying critical data routing paths.

# Evaluation

**Backward**

**Type 1**

**Type 2**

**Partial Extraction**

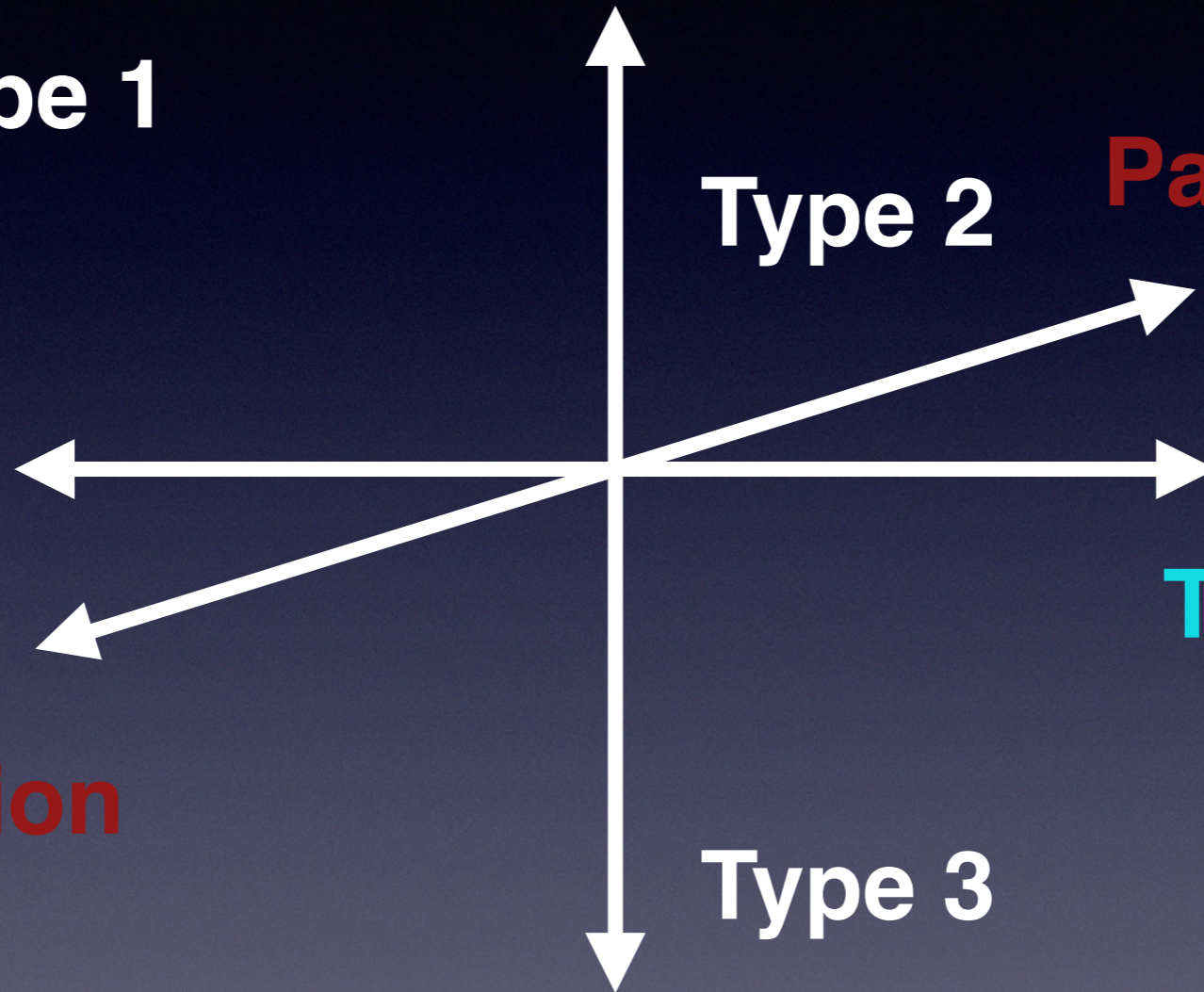
**Sorting**

**Thresholding**

**Full Extraction**

**Type 3**

**Forward**

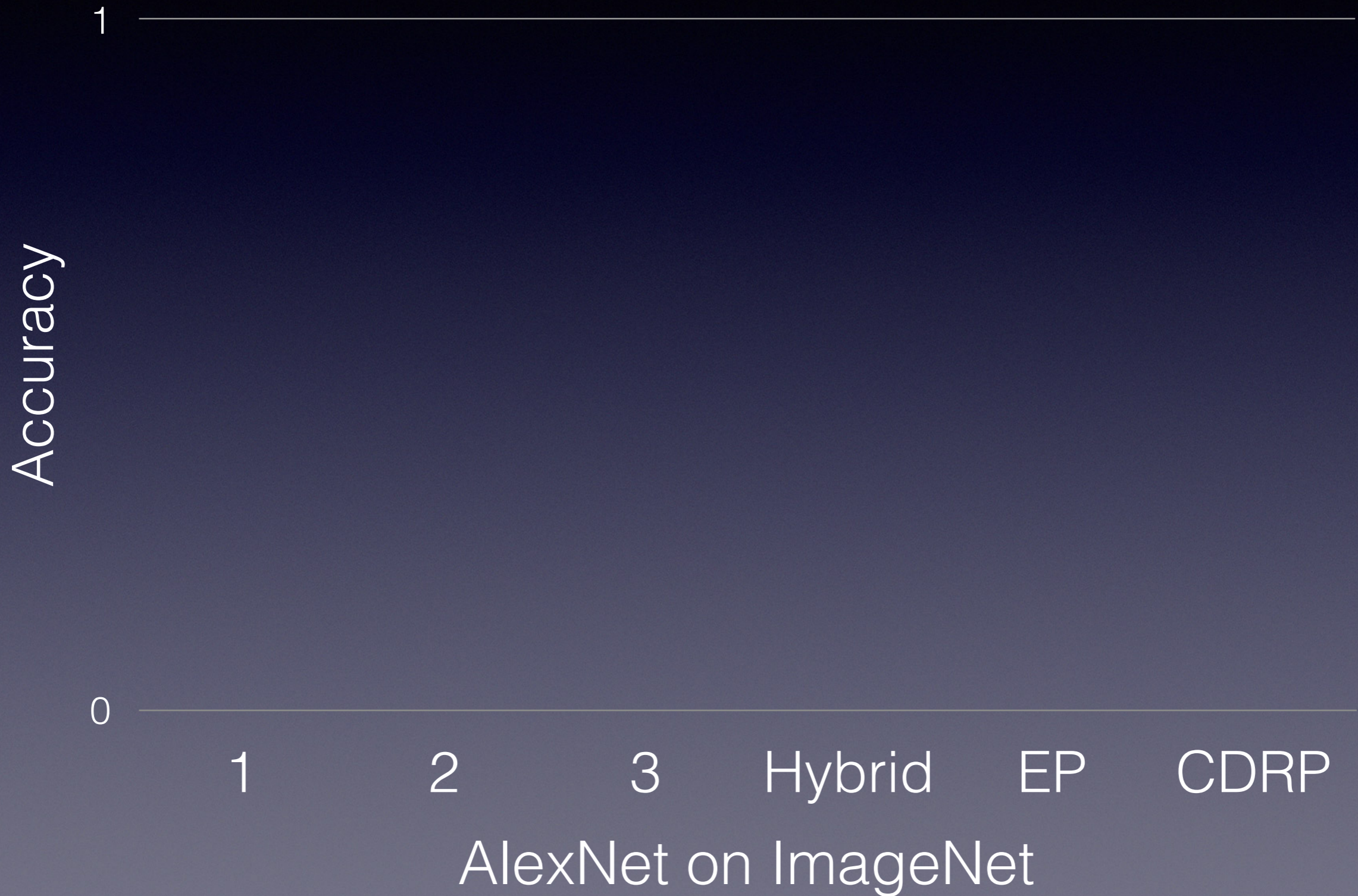




# Hardware Setup

DNN Accelerator	20 x 20
Technology	Silvaco 15nm
On-chip SRAM	1.5MB

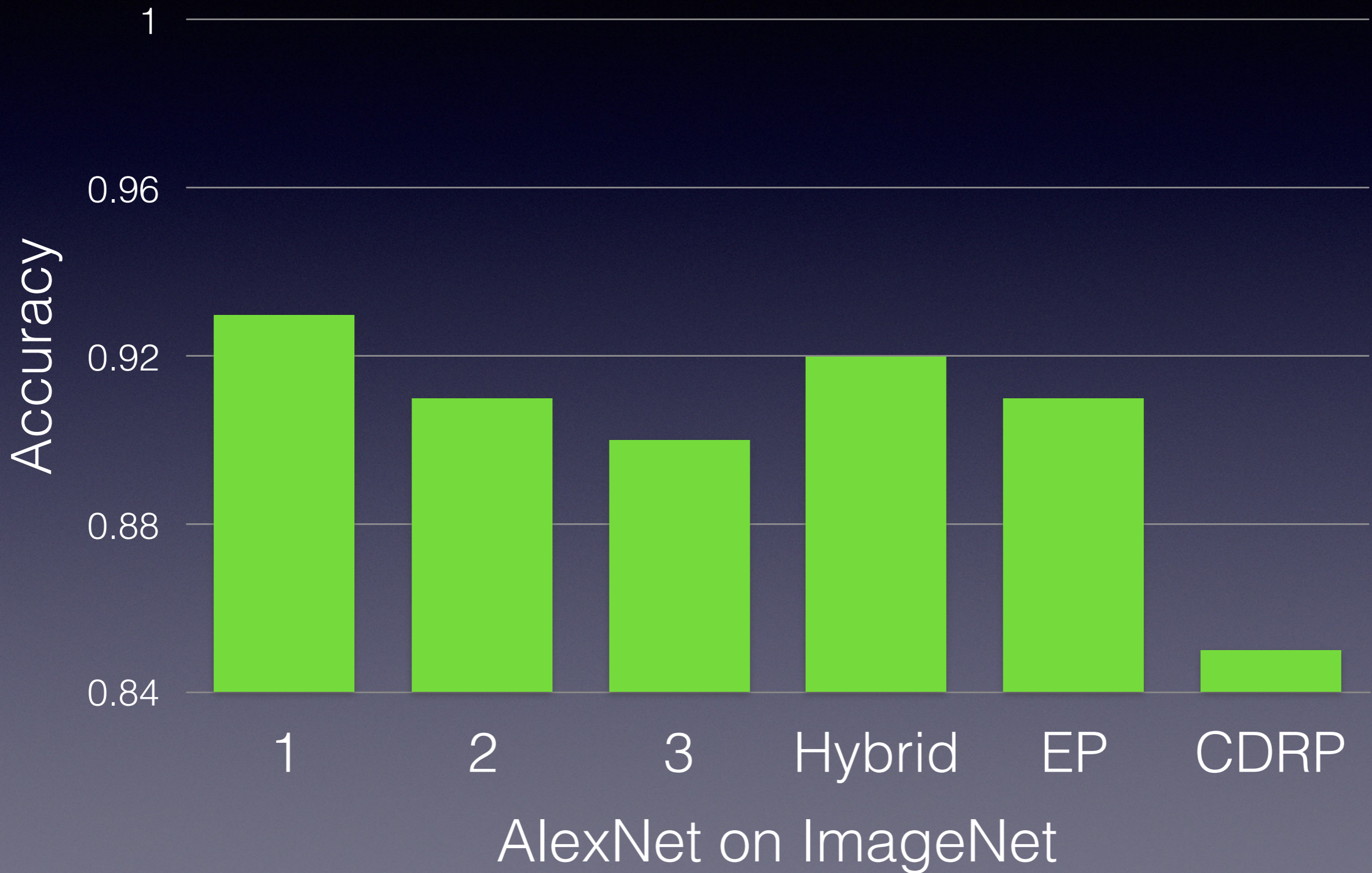
# Evaluation



# Evaluation



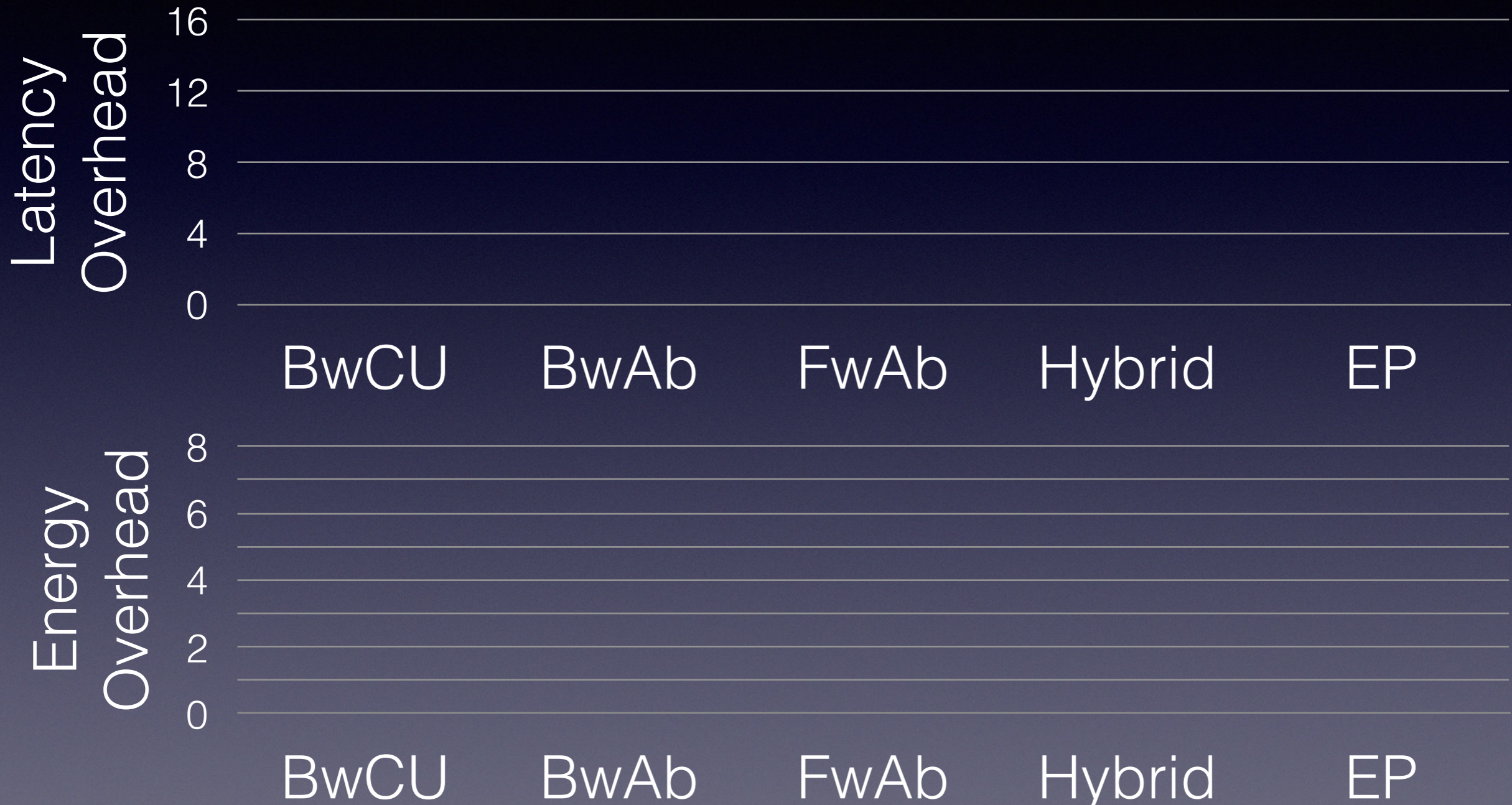
# Evaluation



# Evaluation

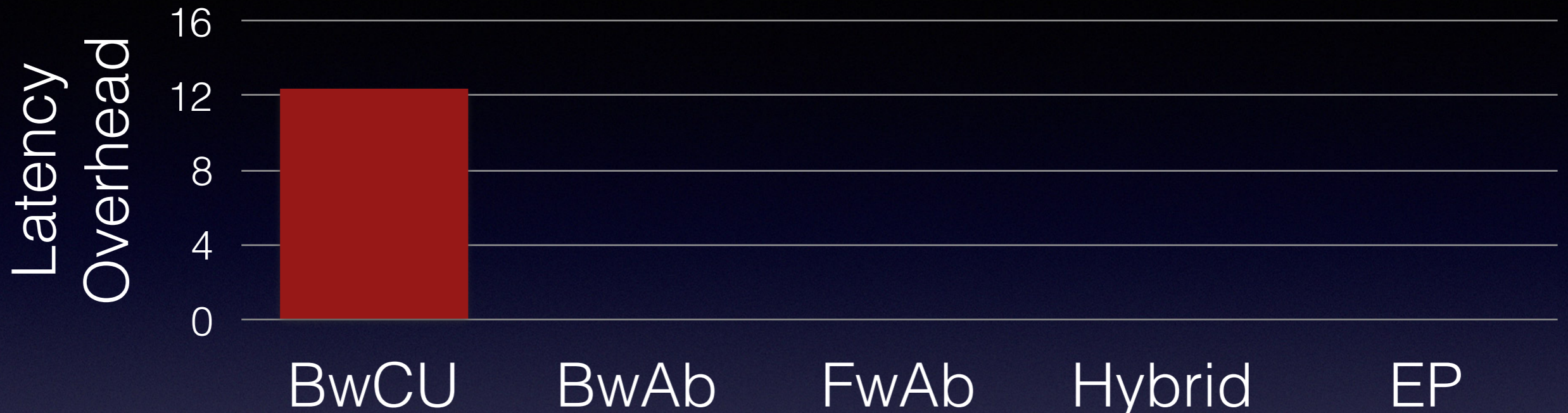


# Evaluation



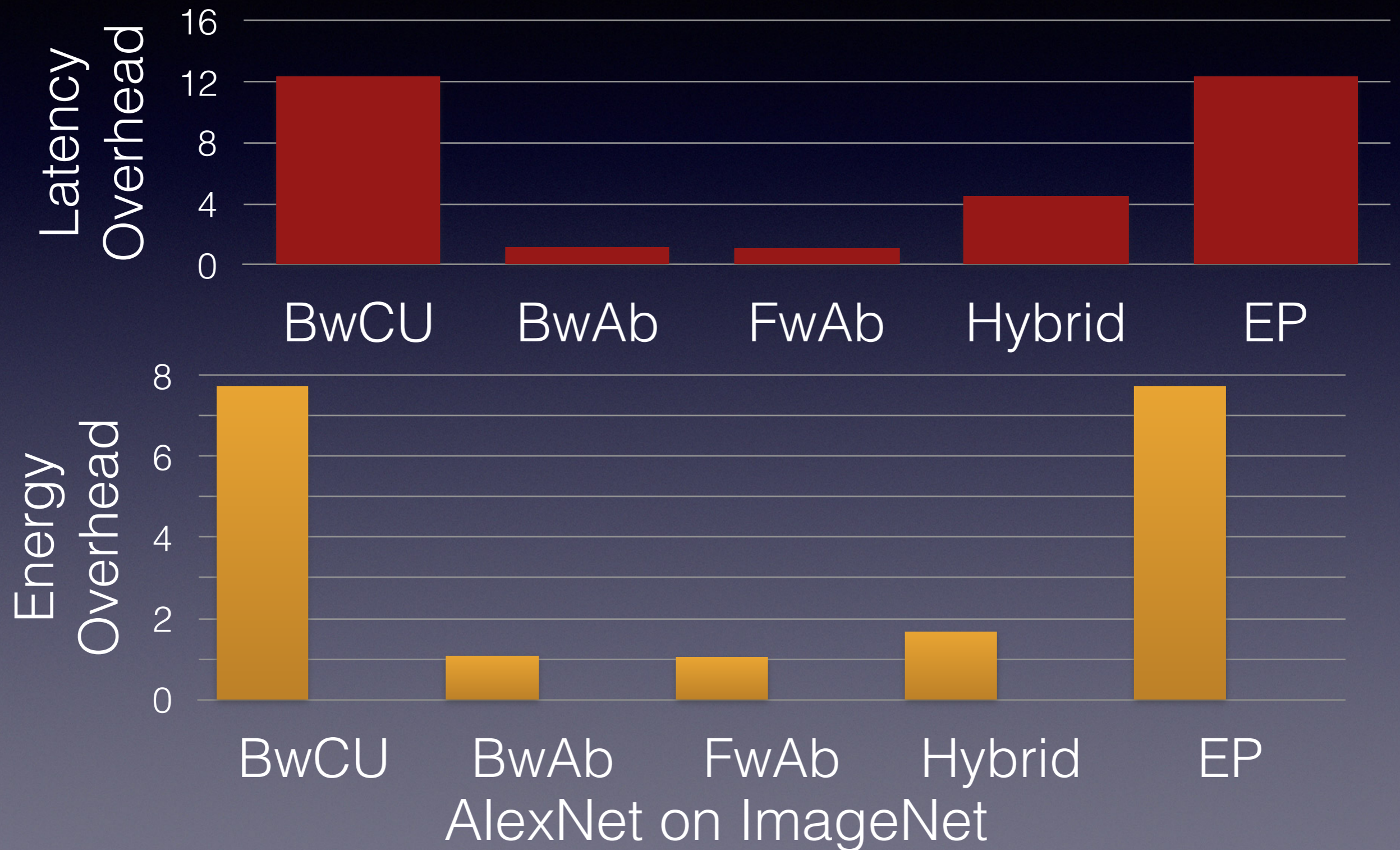
AlexNet on ImageNet

# Evaluation



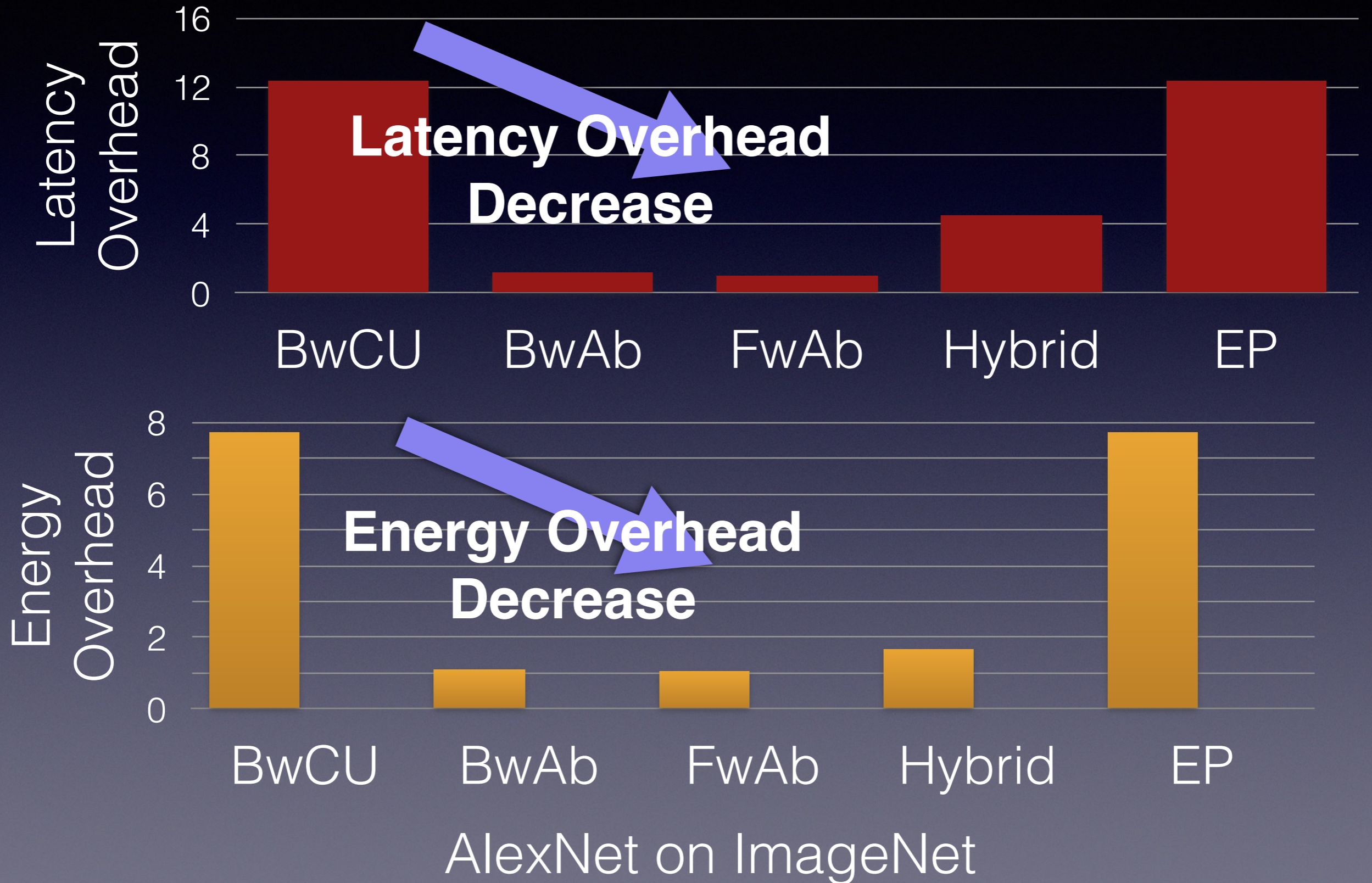
AlexNet on ImageNet

# Evaluation





# Evaluation



# Conclusion

Ptolemy: Accurate, low overhead, adversarial attack detection

- Algorithm Framework
- Compiler Optimization
- Architecture Support

# Collaborators



Yuxian Qiu

Jingwen Leng

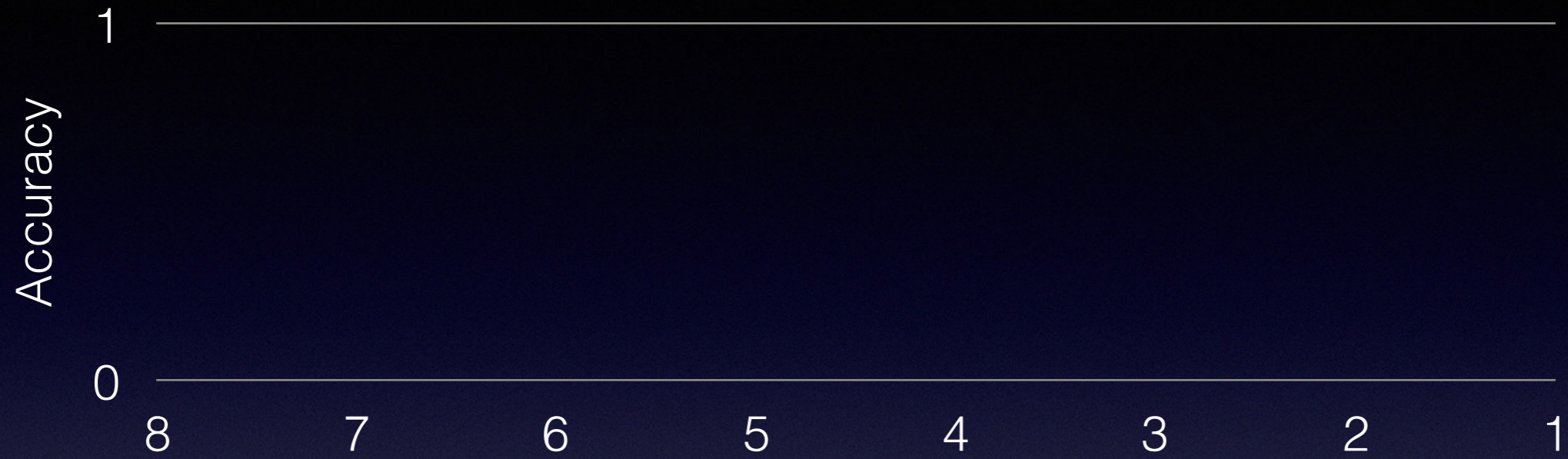
Minyi Guo

Yuhao Zhu

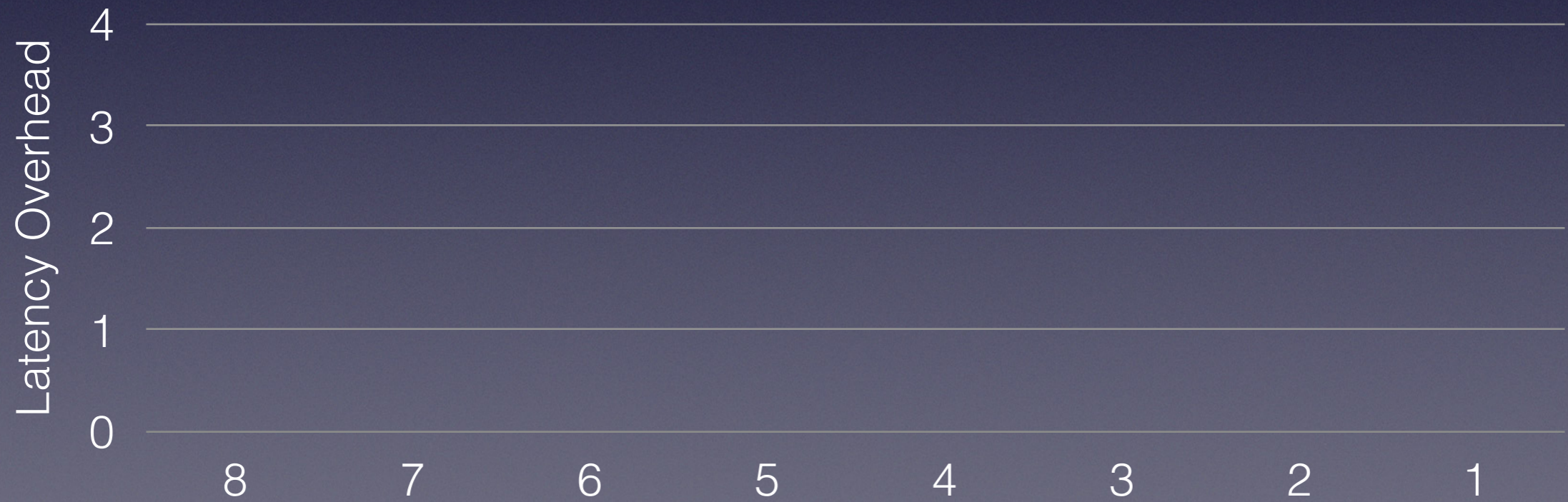
# Questions

<https://github.com/Ptolemy-dl/Ptolemy>

# Evaluation

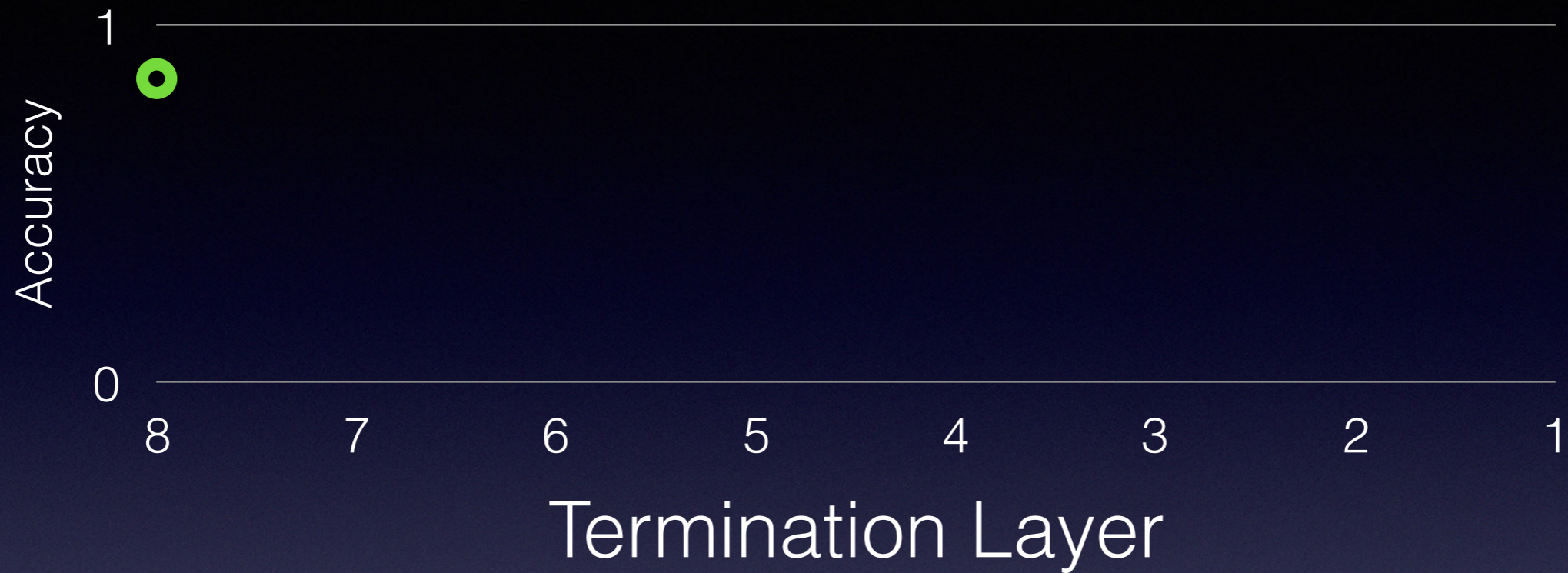


Termination Layer

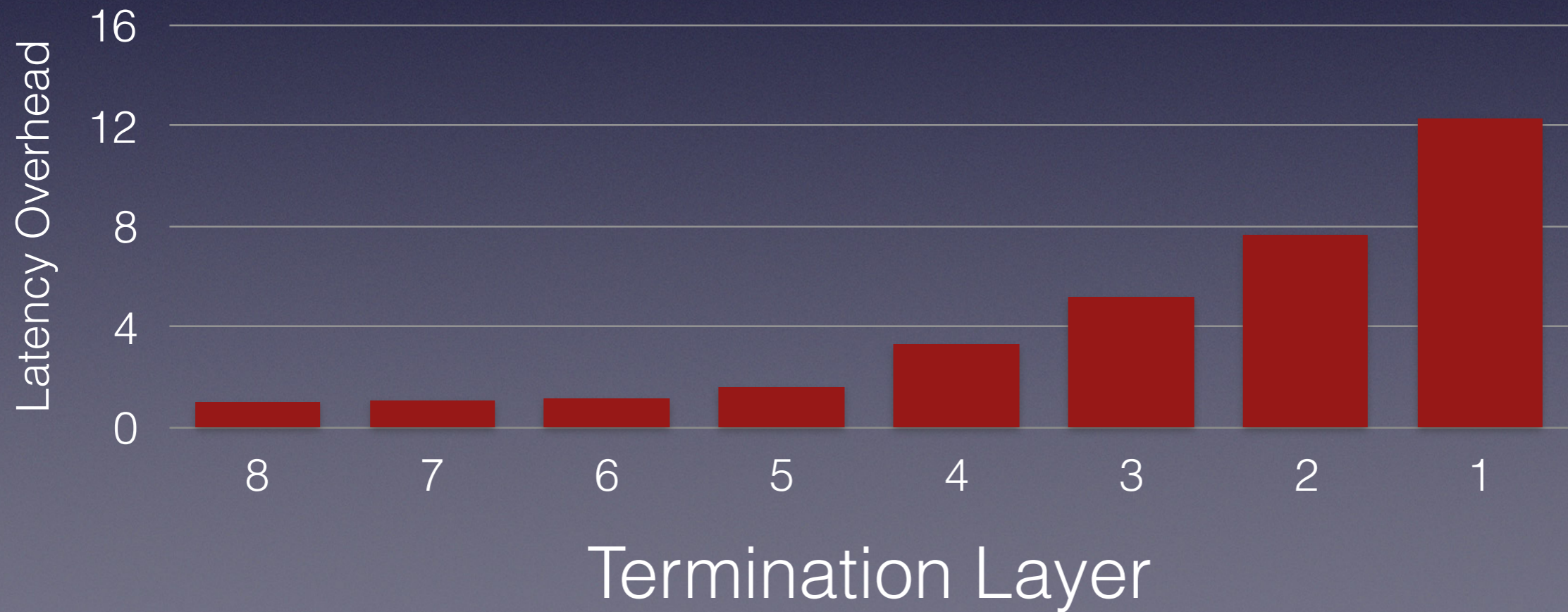
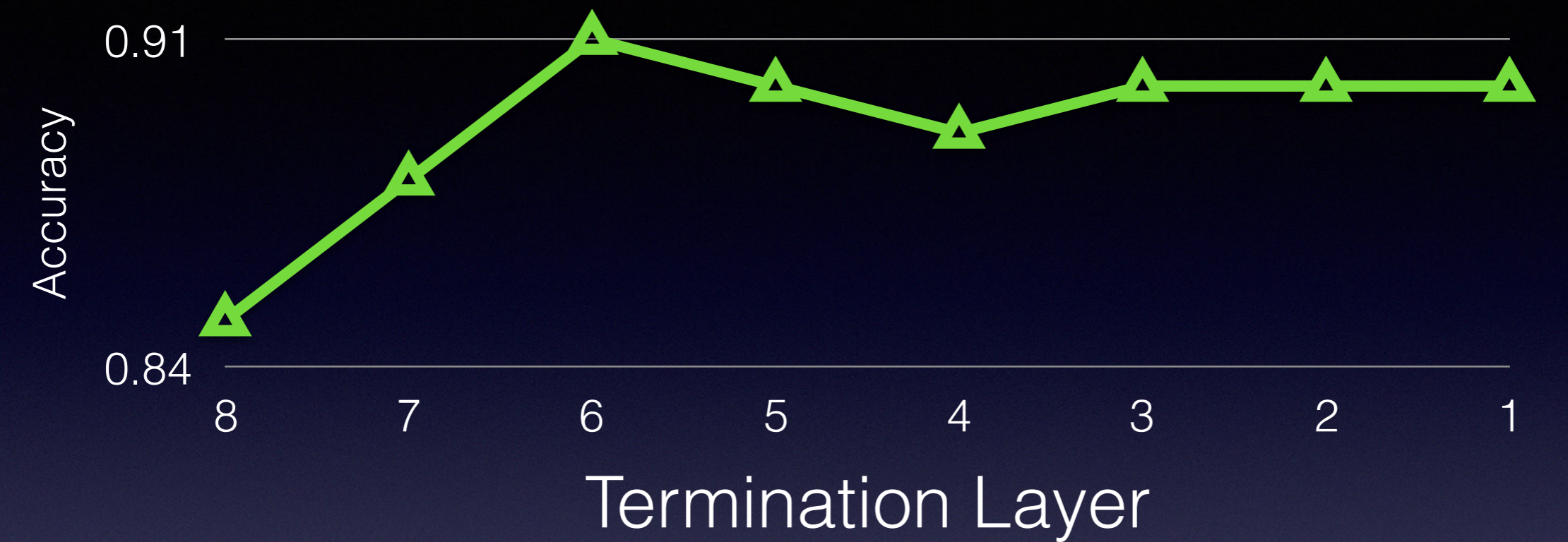


Termination Layer

# Evaluation



# Evaluation



# Backup

```
1 def AdversaryDetection(model, input,  $\theta$ ,  $\phi$ ):
2     output = Inference(model, input)
3     N = model.num_layers
4     // Selective extraction only in the last three layers
5     for L in range(N-3, N):
6         if L != N-1:
7             // Forward extraction using absolute thresholds
8             ImptN[L] = ExtractImptNeurons(1, 1,  $\phi$ , L)
9         else:
10            // Forward extraction using cumulative thresholds
11            ImptN[L] = ExtractImptNeurons(1, 0,  $\theta$ , L)
12            dynPath.concat(GenMask(ImptN[L]))
13    classPath = LoadClassPath(argmax(output))
14    is_adversary = Classify(classPath, dynPath)
15    return is_adversary
```



# Backup

Class	Name	23-20	19-16	15-12	11-8	7-4	3-0
Inference	<b>inf</b>	0000	Input addr.	Weight addr.	Output addr.	Unused	
	<b>infsp</b>	0001	Input addr.	Weight addr.	Output addr.	First partial sum addr.	Unused
	<b>csps</b>	0010	Output neuron ID	Layer ID	First partial sum addr.	Unused	
Path Construction	<b>sort</b>	0011	Unsorted seq. start addr.	Seq. length	Sorted seq. start addr.	Unused	
	<b>acum</b>	0100	Input addr.	Output addr.	Cumulative threshold	Unused	
	<b>genmasks</b>	0101	Input addr.	Output addr.	Unused		
	<b>findneuron</b>	0110	Layer ID	Neuron position	Target neuron addr.	Unused	
	<b>findrf</b>	0111	Neuron addr.	Receptive field addr.	Unused		
Classification	<b>cls</b>	1000	Class path addr.	Activation path addr.	Result	Unused	
Others	Omitted for simplicity ( <b>mov</b> , <b>dec</b> , <b>jne</b> , etc.)						

# Backup

; inference & extraction of a model

```
| for j = 1 to L {  
|   inf(j)  
|   <extraction for j>  
| }
```



```
| inf(1)  
| for j = 1 to L-1 {  
|   inf(j+1)  
|   <extraction for j>  
| }  
| <extraction for L>
```

(a) Overlapping inference with extraction across layers in forward extraction.  $L$  is the total number of DNN layers.

; extraction within a layer

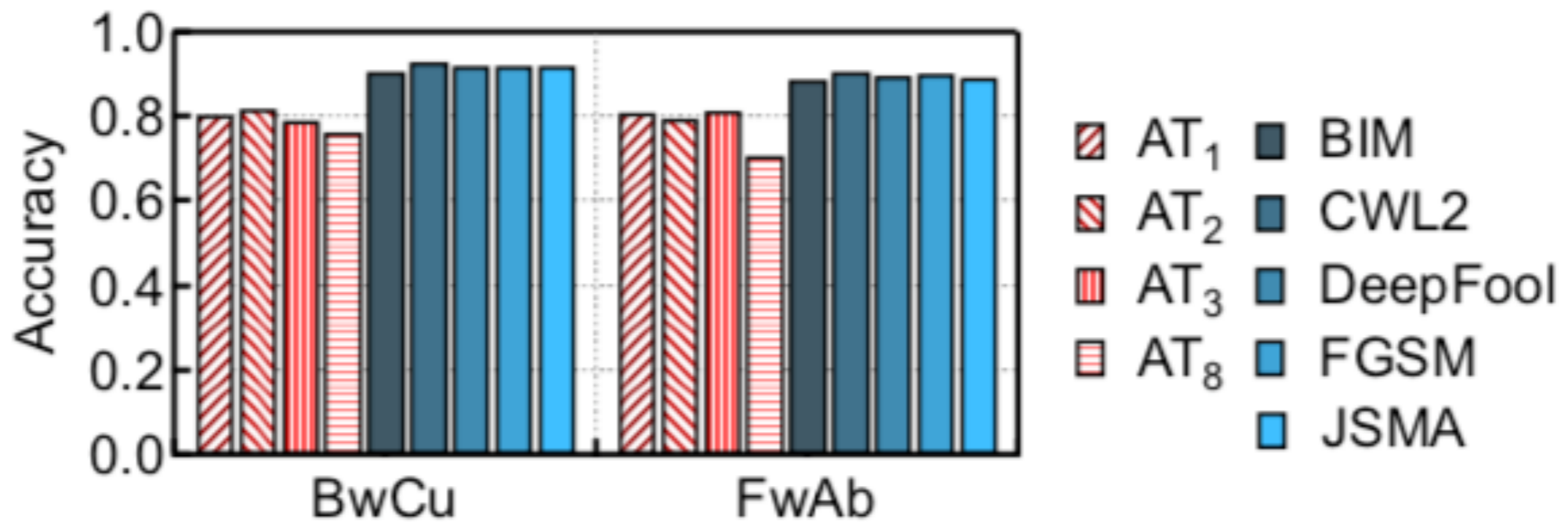
```
| for i = 1 to N {  
|   sort(i)  
|   acum(i)  
| }
```



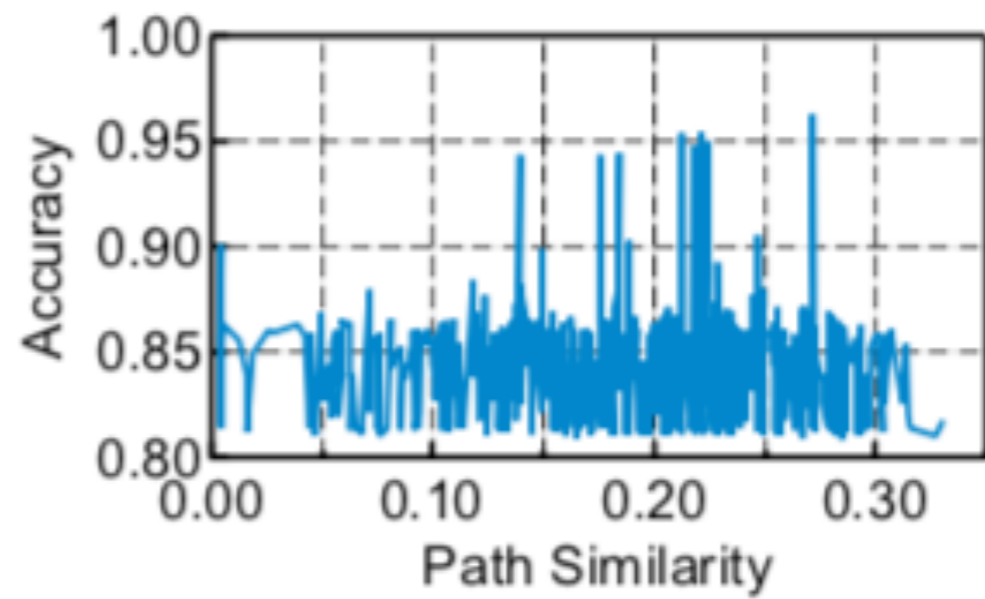
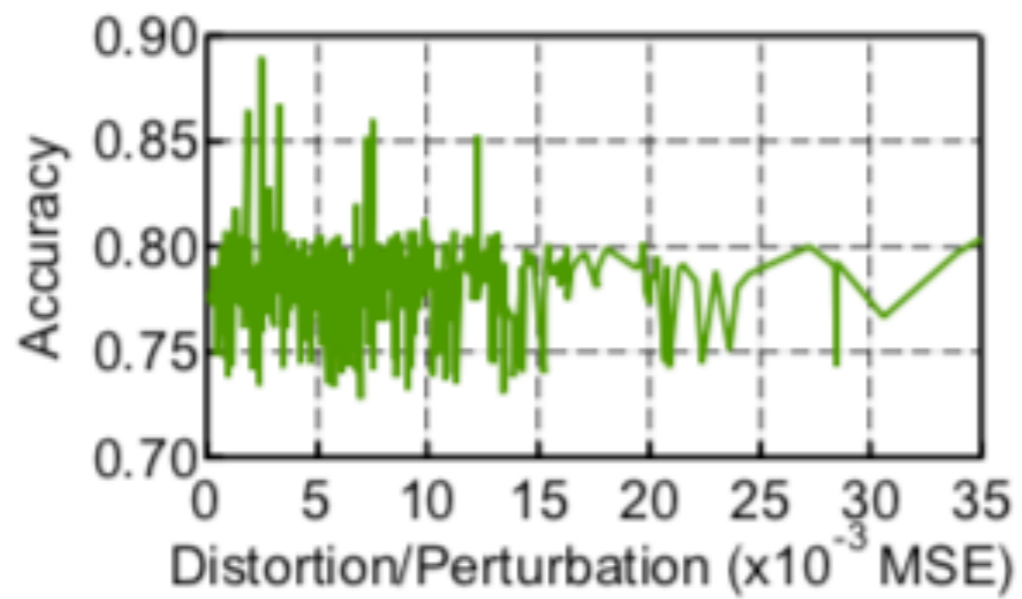
```
| sort(1)  
| for i = 1 to N-1 {  
|   sort(i+1)  
|   acum(i)  
| }  
| acum(N)
```

(b) Neuron-level pipelining in important neuron extraction.  $N$  denotes the number of important neurons in the current layer's output.

# Backup



# Backup



Backup