

LightFF: Lightweight Inference for Forward-Forward Algorithm

Amin Aminifar*, **Baichuan Huang**[†], **Azra Abtahi**[†], **Amir Aminifar**[†]

* Institute of Computer Engineering, Heidelberg University, Germany

† Department of Electrical and Information Technology, Lund University, Sweden

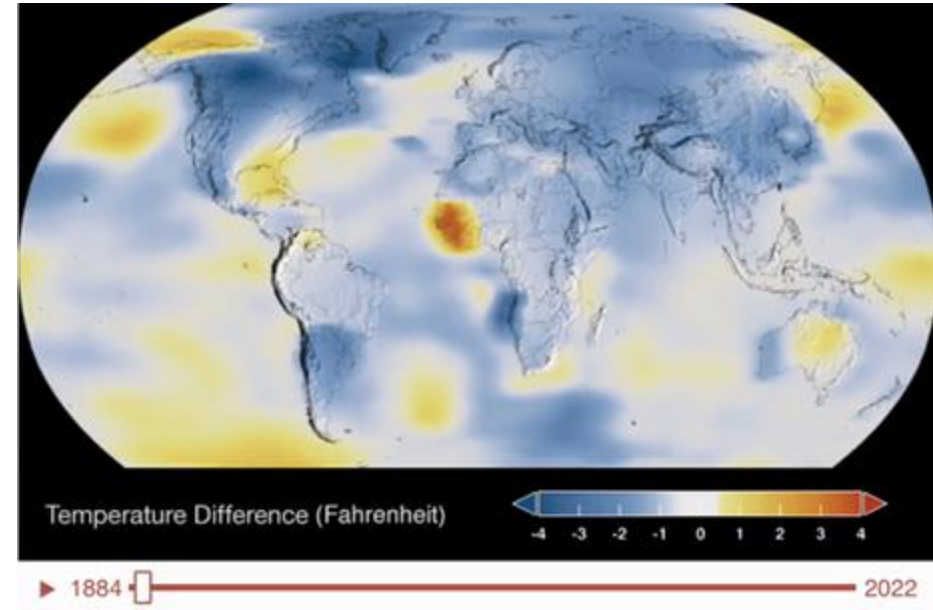
This research has been partially supported by the Swedish Wallenberg AI, Autonomous Systems and Software Program (WASP), the Swedish Research Council (VR), Swedish Foundation for Strategic Research (SSF), the ELLIIT Strategic Research Environment, and the European Union (EU).



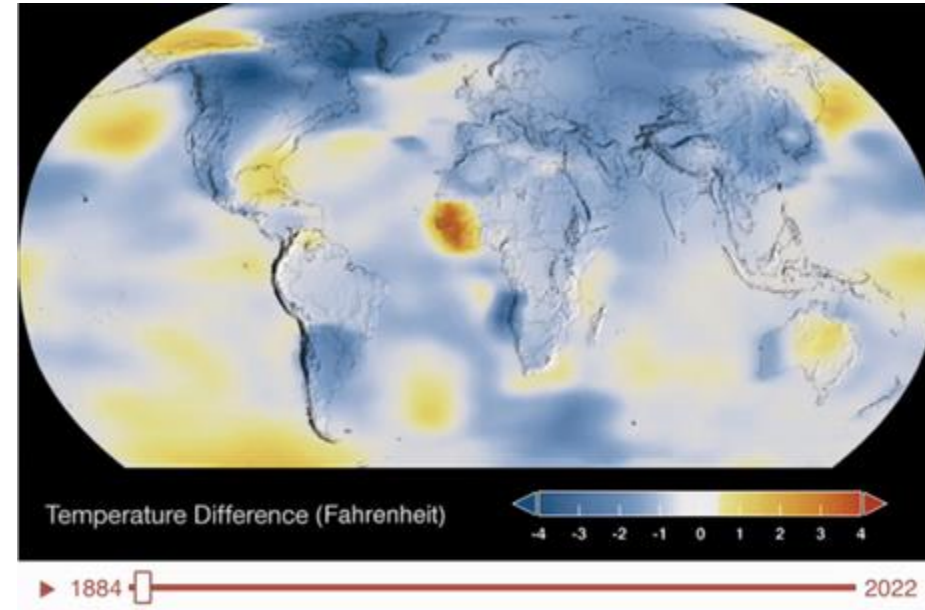
LUND
UNIVERSITY

Introduction and Background

Global Warming

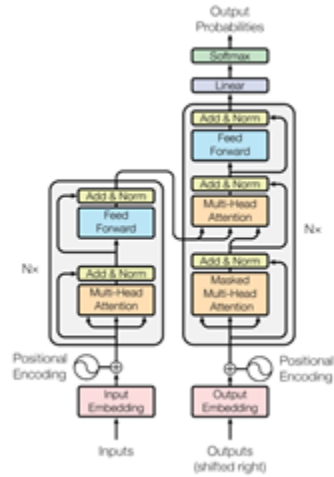


Global Warming



Europe: an average rise of 2.3°C compared to pre-industrial levels
 1°C **higher than** the global average.

Environmental Impact of Deep Learning



Training Transformer (Strubell E. 2020)



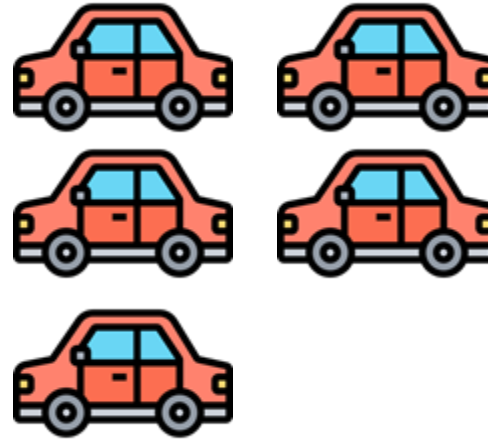
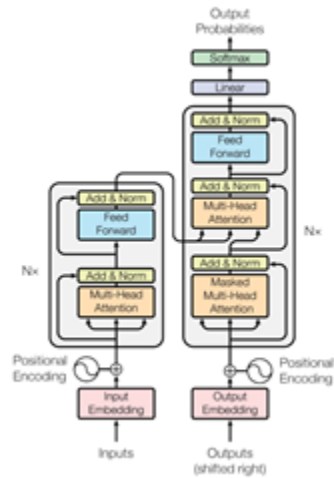
626,155 lbs

Strubell E, et al. Energy and policy considerations for modern deep learning research. AAAI, 2020.

Vaswani A. Attention is all you need. NeurIPS, 2017.

<https://www.forbes.com/sites/robtnews/2020/06/17/deep-learnings-climate-change-problem/>

Environmental Impact of Deep Learning



Training Transformer (Strubell E. 2020)



626,155 lbs

=

Total Lifetime of a Car

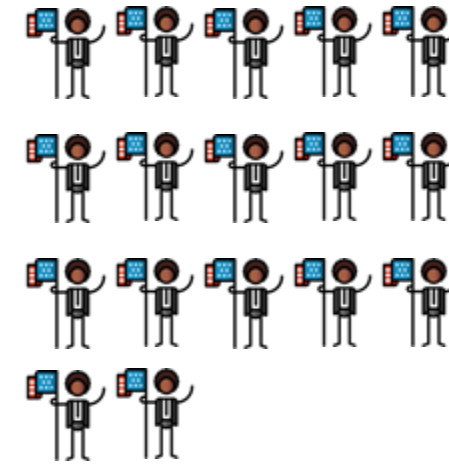
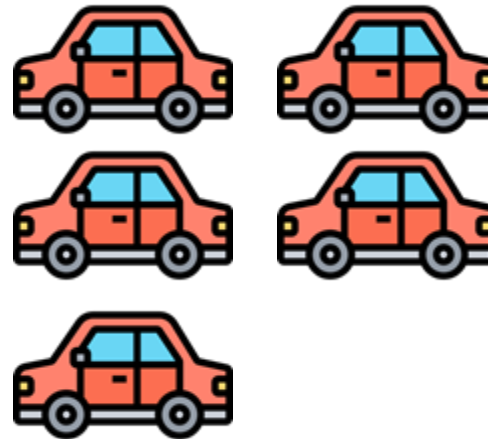
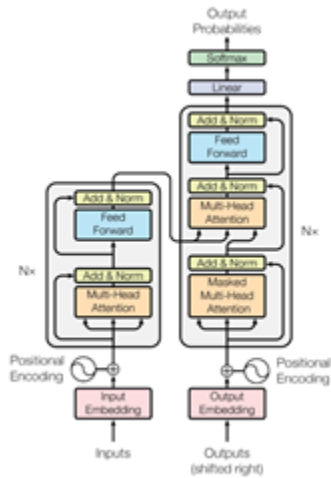
5 × 126,000 lbs

Strubell E, et al. Energy and policy considerations for modern deep learning research. AAAI, 2020.

Vaswani A. Attention is all you need. NeurIPS, 2017.

<https://www.forbes.com/sites/robtnews/2020/06/17/deep-learnings-climate-change-problem/>

Environmental Impact of Deep Learning



Training Transformer (Strubell E. 2020)



626,155 lbs

=

5 × 126,000 lbs

=

17 × 36,156 lbs

Total Lifetime of a Car

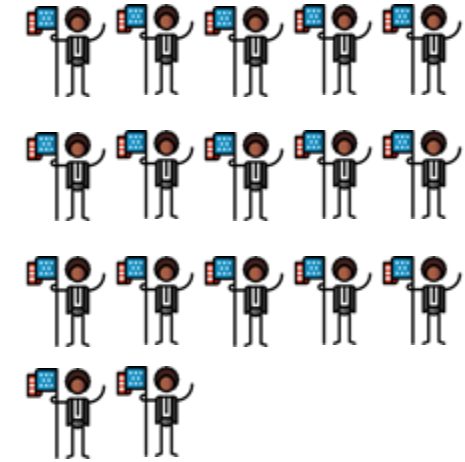
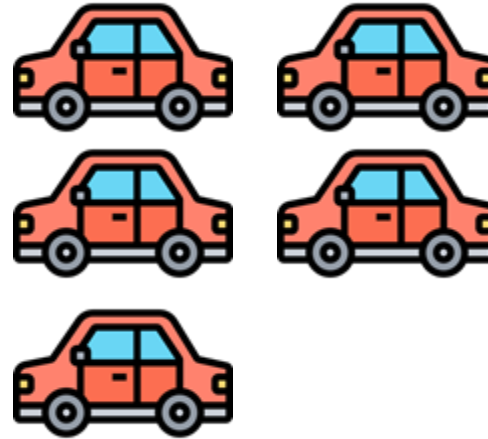
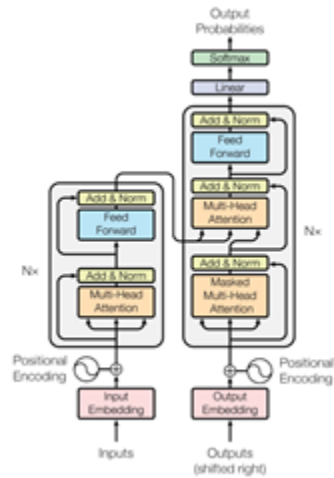
Average American in a Year

Strubell E, et al. Energy and policy considerations for modern deep learning research. AAAI, 2020.

Vaswani A. Attention is all you need. NeurIPS, 2017.

<https://www.forbes.com/sites/robtnews/2020/06/17/deep-learning-climate-change-problem/>

Environmental Impact of Deep Learning



Training Transformer (Strubell E. 2020)

Total Lifetime of a Car

Average American in a Year



626,155 lbs

=

5 × 126,000 lbs

=

17 × 36,156 lbs

The computational resources needed to produce a best-in-class AI model has on average **doubled every 3.4 months.**

Strubell E, et al. Energy and policy considerations for modern deep learning research. AAAI, 2020.

Vaswani A. Attention is all you need. NeurIPS, 2017.

<https://www.forbes.com/sites/robtowes/2020/06/17/deep-learnings-climate-change-problem/>

Energy Consumption of Training



GPT-3



GPT-4

D. Patterson, et al. Carbon emissions and large neural network training, 2021.

<https://tinymml.substack.com/p/the-carbon-impact-of-large-language>

Data sources: U.S. Energy Information Administration, Electric Power Research Institute (EPRI)

Energy Consumption of Training



GPT-3



GPT-4



1,216,950 lbs

×13

15,238,333 lbs

D. Patterson, et al. Carbon emissions and large neural network training, 2021.

<https://tinymml.substack.com/p/the-carbon-impact-of-large-language>

Data sources: U.S. Energy Information Administration, Electric Power Research Institute (EPRI)

Energy Consumption of Training



GPT-3



GPT-4



1,216,950 lbs

×13

15,238,333 lbs



1,287 Megawatt-Hour

×48

62,318 Megawatt-Hour

D. Patterson, et al. Carbon emissions and large neural network training, 2021.

<https://tinymml.substack.com/p/the-carbon-impact-of-large-language>

Data sources: U.S. Energy Information Administration, Electric Power Research Institute (EPRI)

Biologically Plausible Alternatives

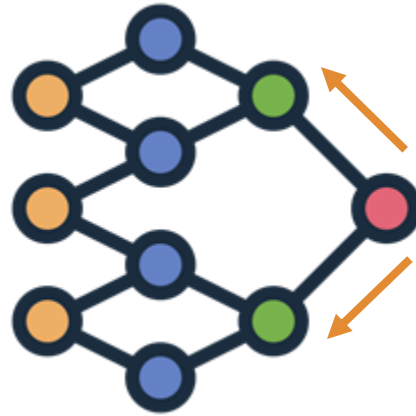


Human Brain
(~20 Watts)

Biologically Plausible Alternatives



Human Brain
(~20 Watts)

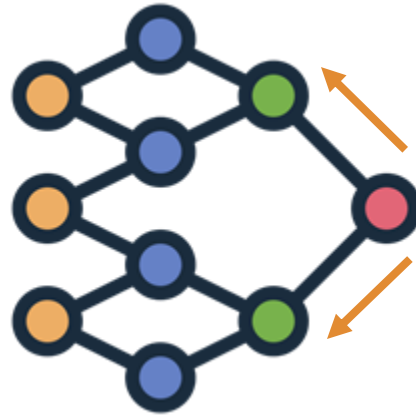


Back-Propagation
(Bio-**Implausible**)

Biologically Plausible Alternatives



Human Brain
(~20 Watts)



Back-Propagation
(Bio-Implausible)



Forward-Forward Algorithm
(Bio-Plausible)

Energy Consumption of Inference



~60%



80%-90%



~90%

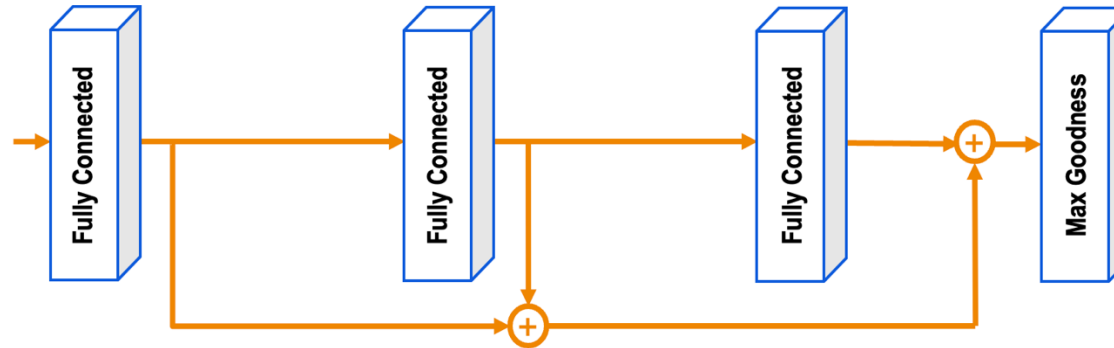
The **percentage** of inference in the total ML energy consumed.



LUND
UNIVERSITY

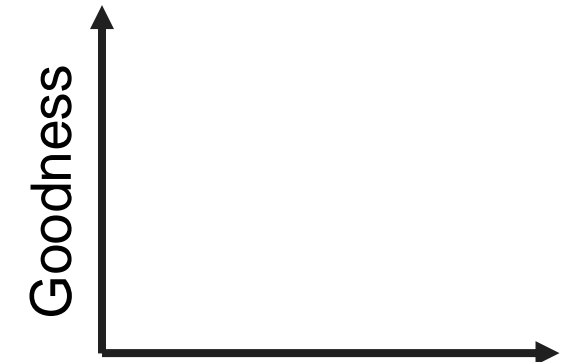
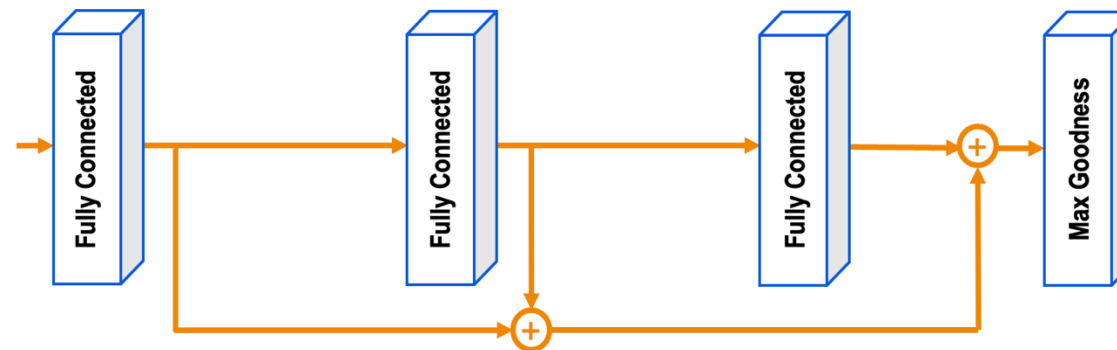
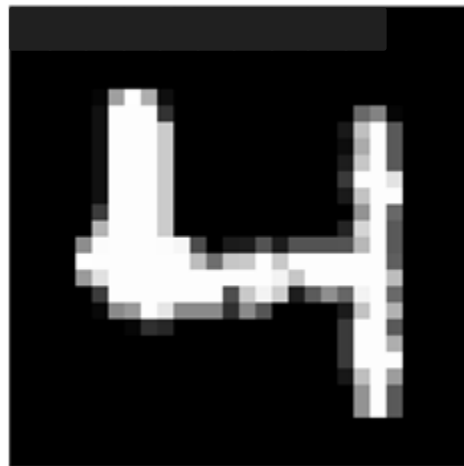
Lightweight Inference for Forward-Forward Algorithm

Glance at Hinton's FF Inference



The Goodness is the sum of the square of the activity of each hidden neuron.

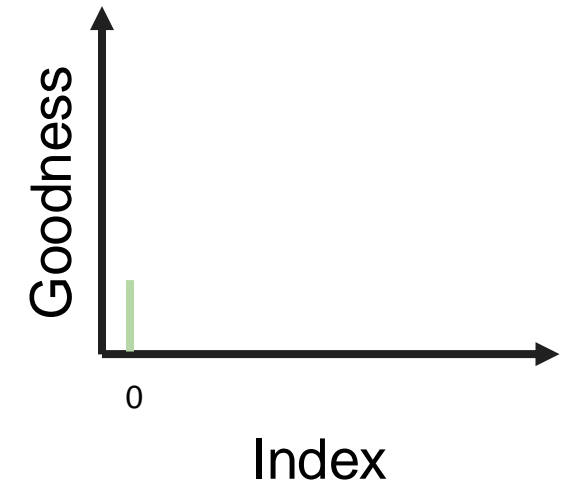
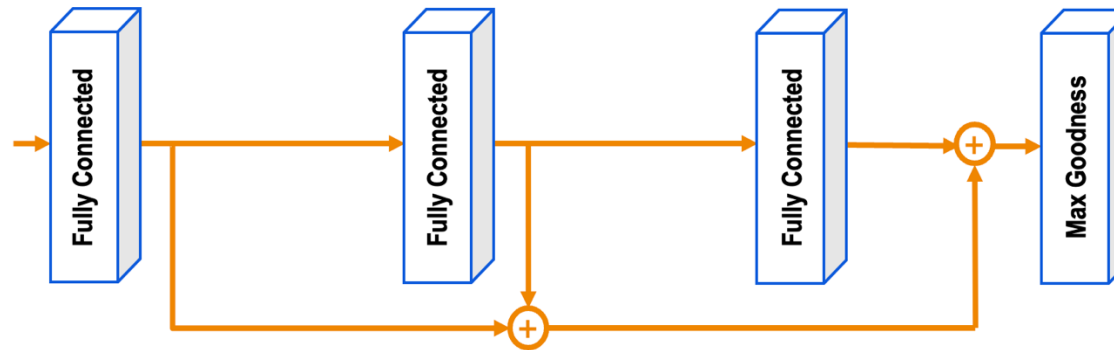
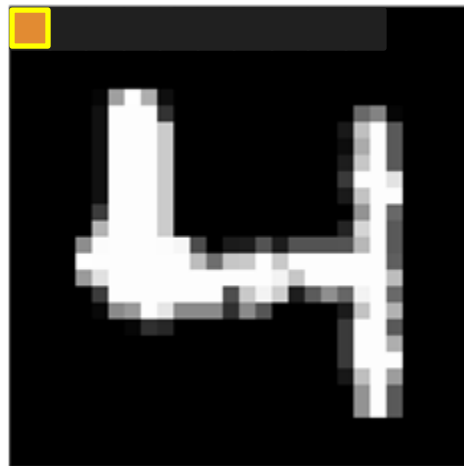
Glance at Hinton's FF Inference



The Goodness is the sum of the square of the activity of each hidden neuron.

Glance at Hinton's FF Inference

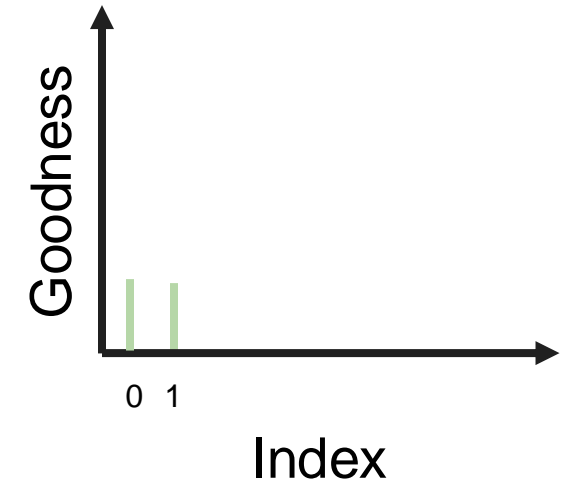
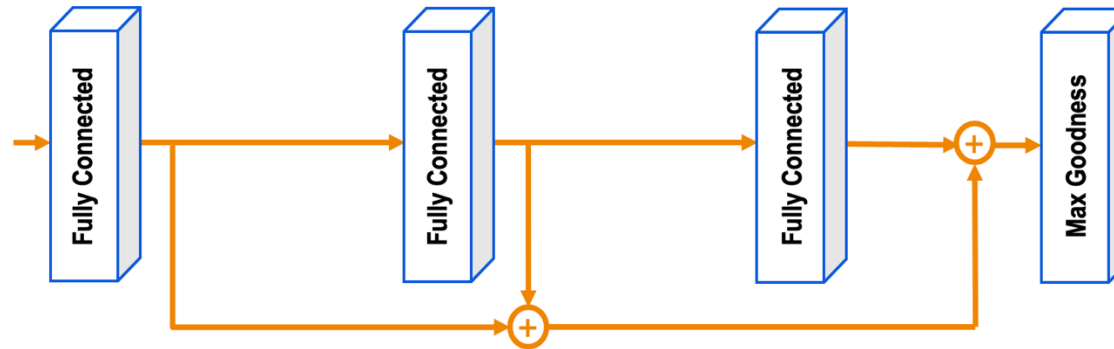
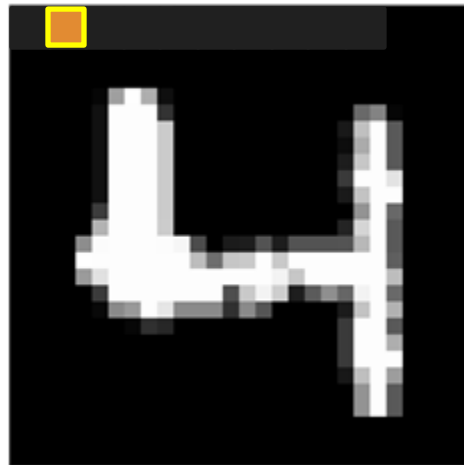
Index: 0



The Goodness is the sum of the square of the activity of each hidden neuron.

Glance at Hinton's FF Inference

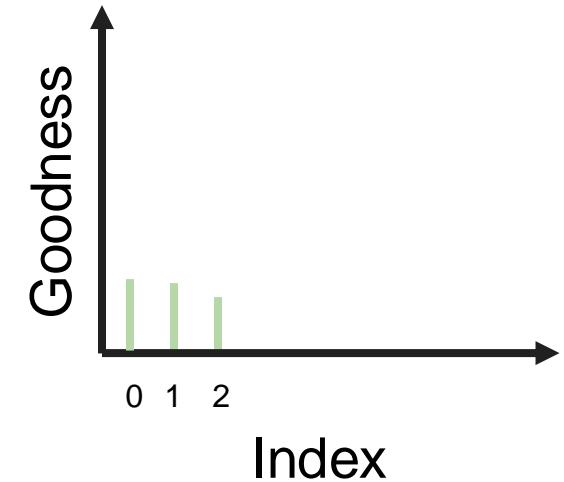
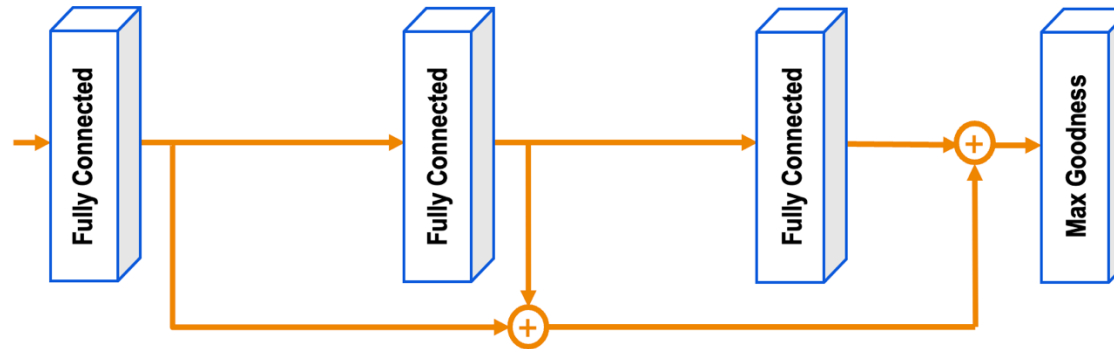
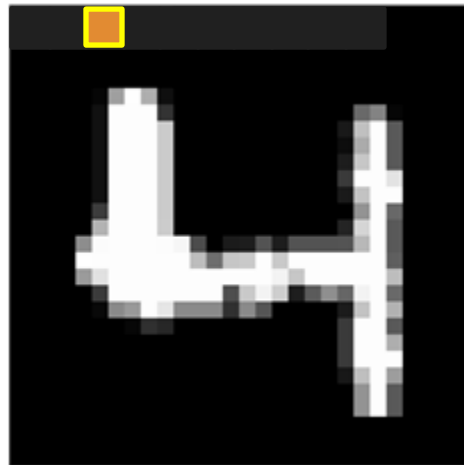
Index: 1



The Goodness is the sum of the square of the activity of each hidden neuron.

Glance at Hinton's FF Inference

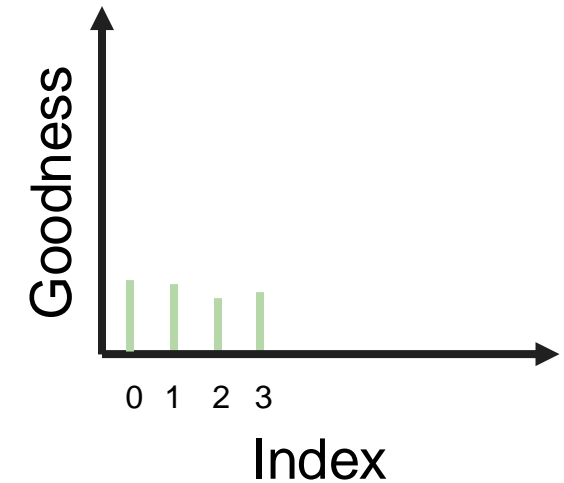
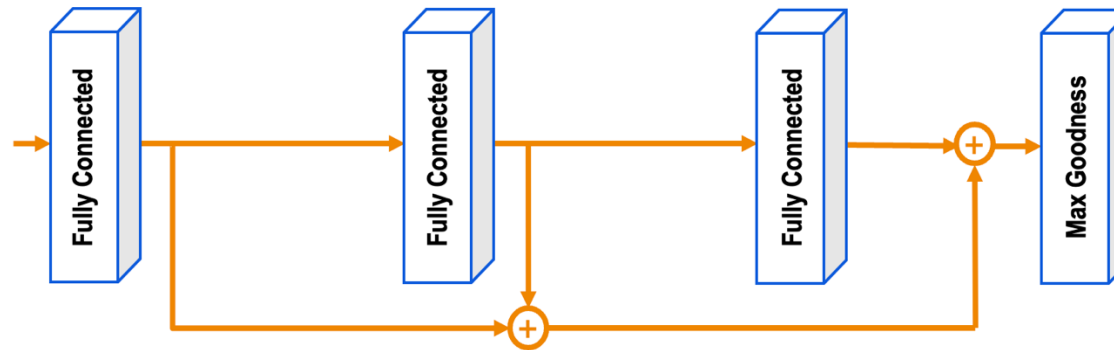
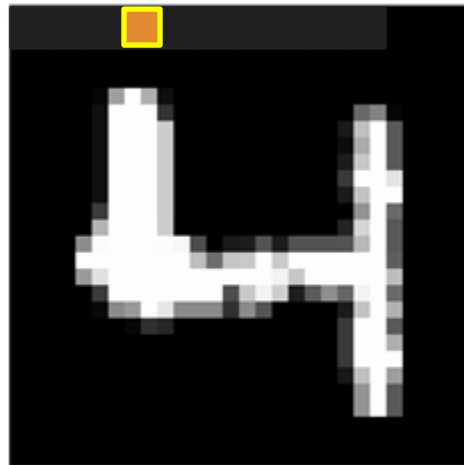
Index: 2



The Goodness is the sum of the square of the activity of each hidden neuron.

Glance at Hinton's FF Inference

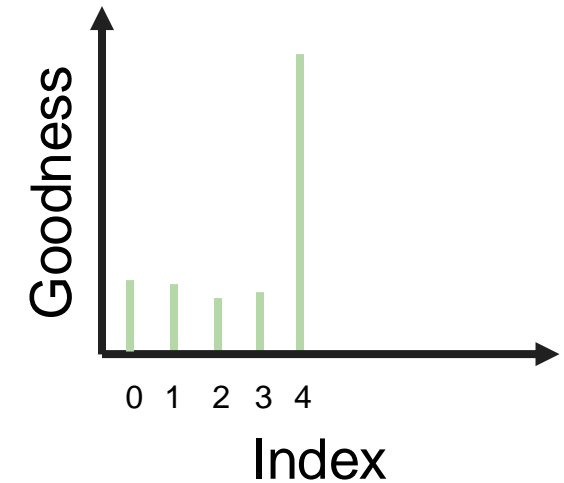
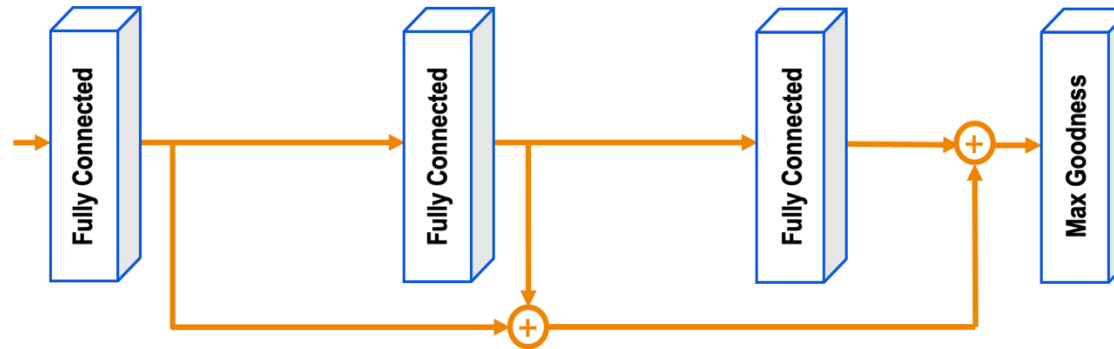
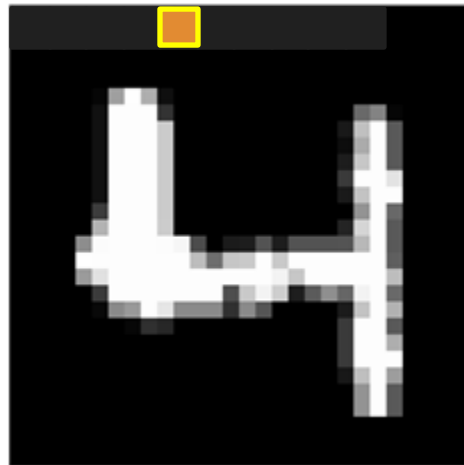
Index: 3



The Goodness is the sum of the square of the activity of each hidden neuron.

Glance at Hinton's FF Inference

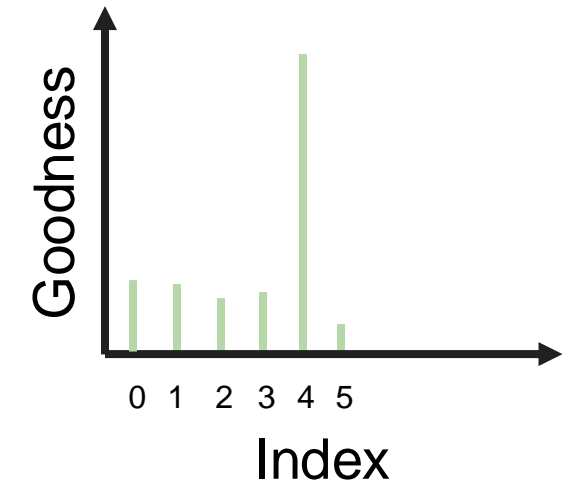
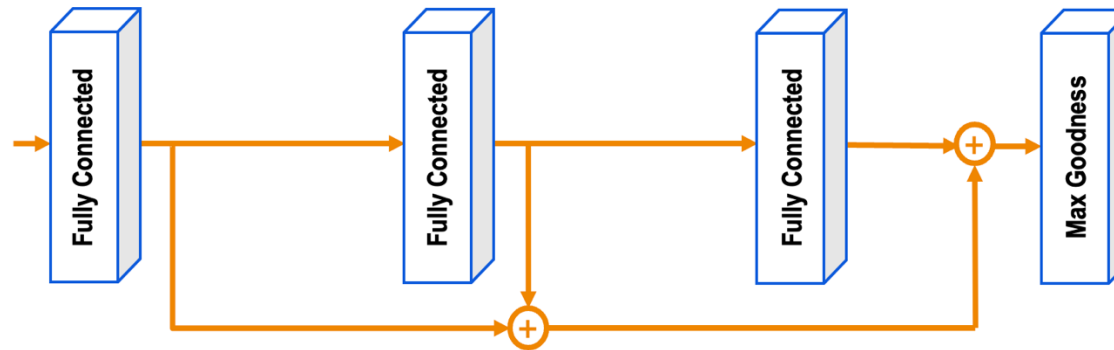
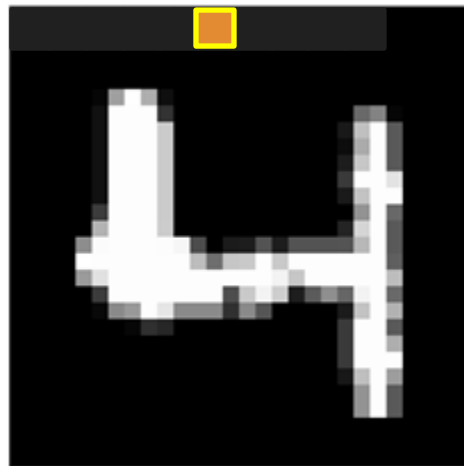
Index: 4



The Goodness is the sum of the square of the activity of each hidden neuron.

Glance at Hinton's FF Inference

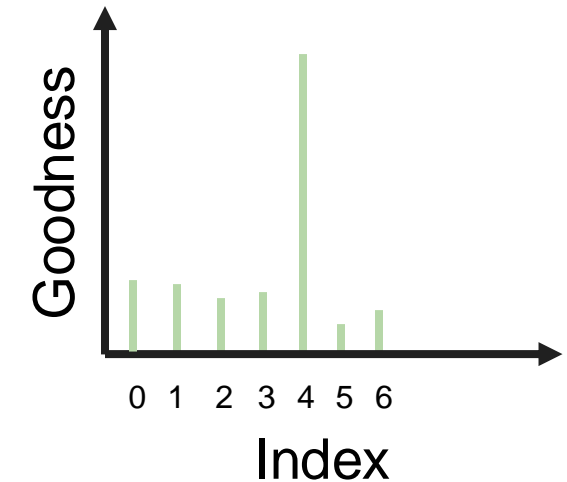
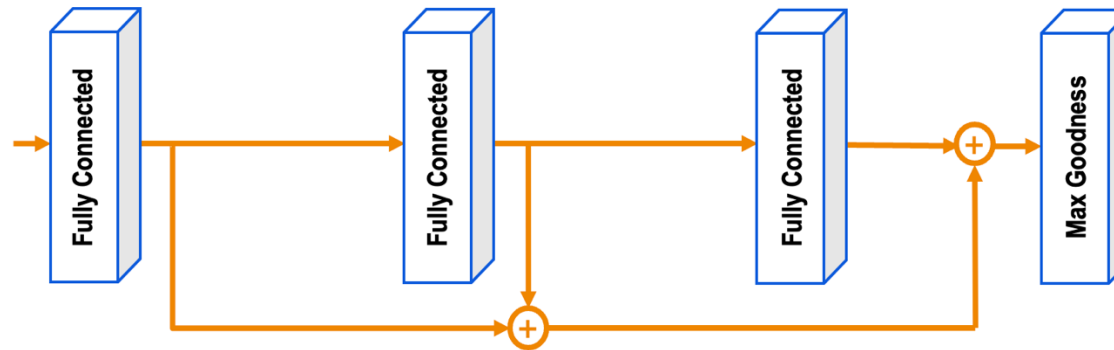
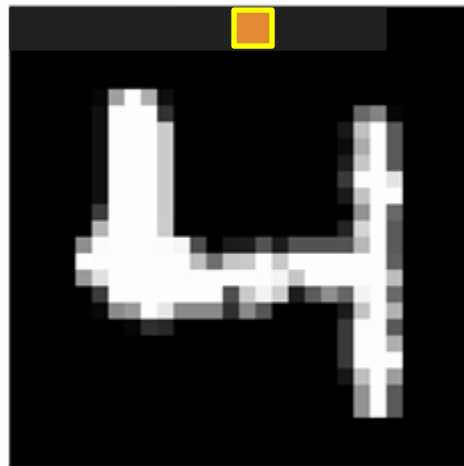
Index: 5



The Goodness is the sum of the square of the activity of each hidden neuron.

Glance at Hinton's FF Inference

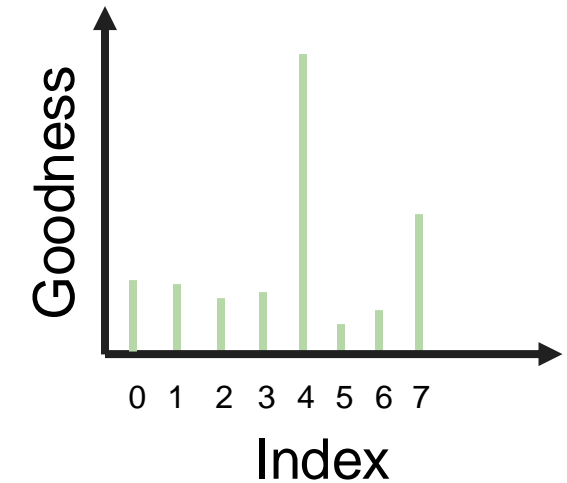
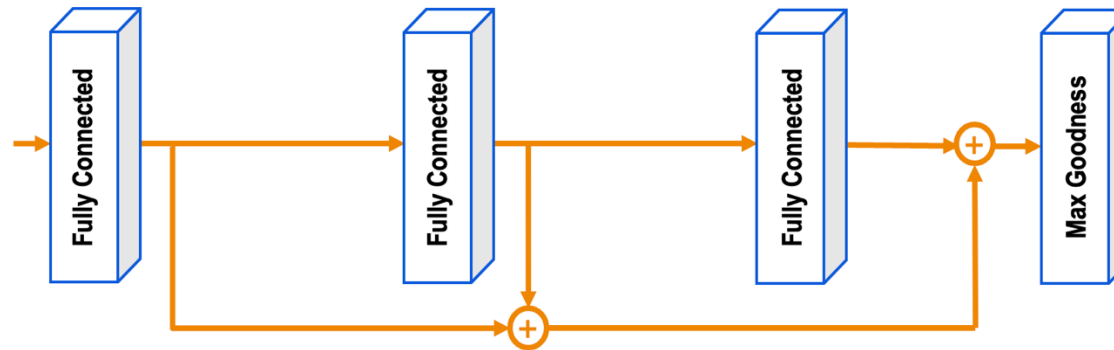
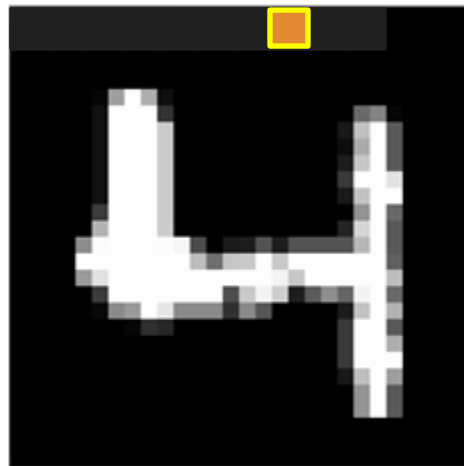
Index: 6



The Goodness is the sum of the square of the activity of each hidden neuron.

Glance at Hinton's FF Inference

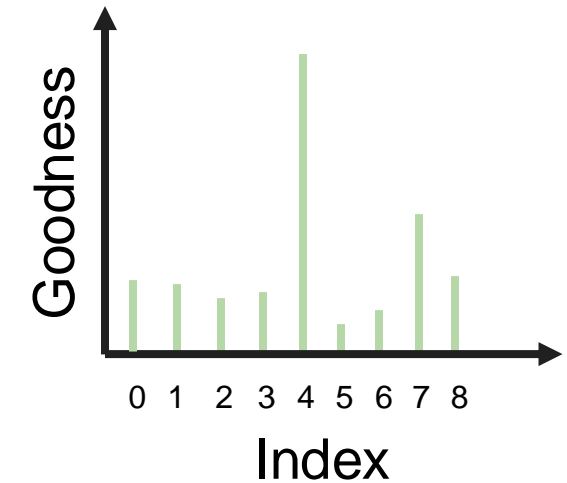
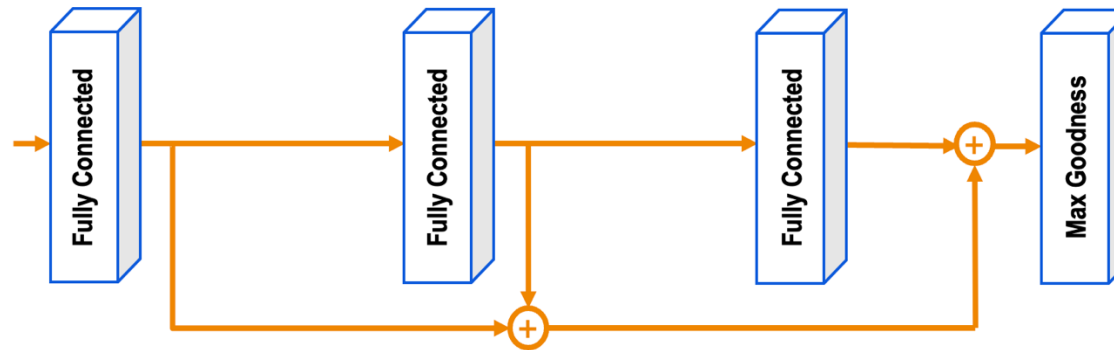
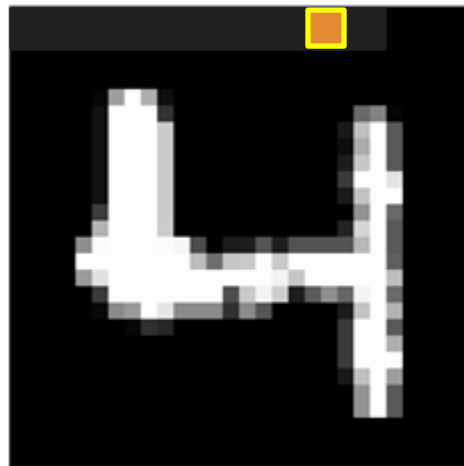
Index: 7



The Goodness is the sum of the square of the activity of each hidden neuron.

Glance at Hinton's FF Inference

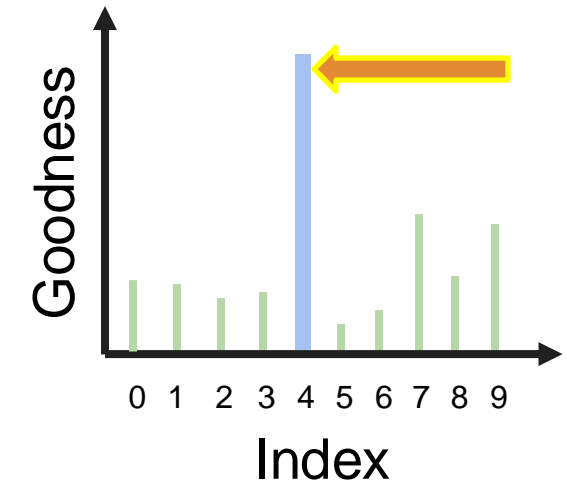
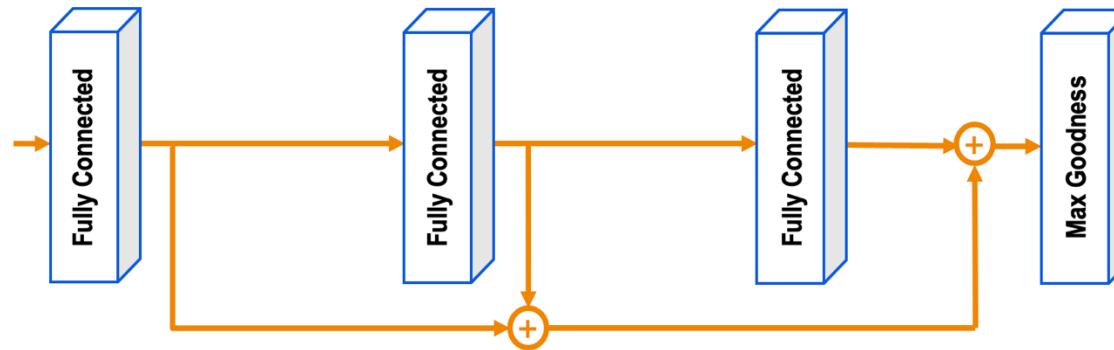
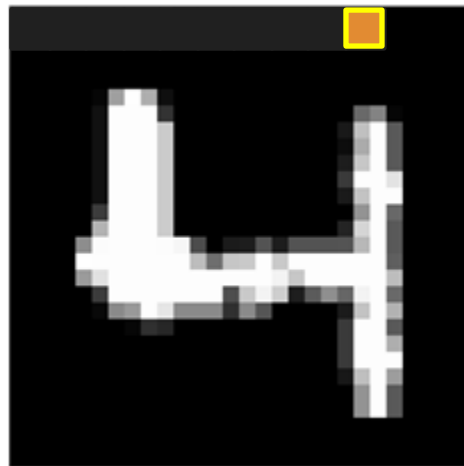
Index: 8



The Goodness is the sum of the square of the activity of each hidden neuron.

Glance at Hinton's FF Inference

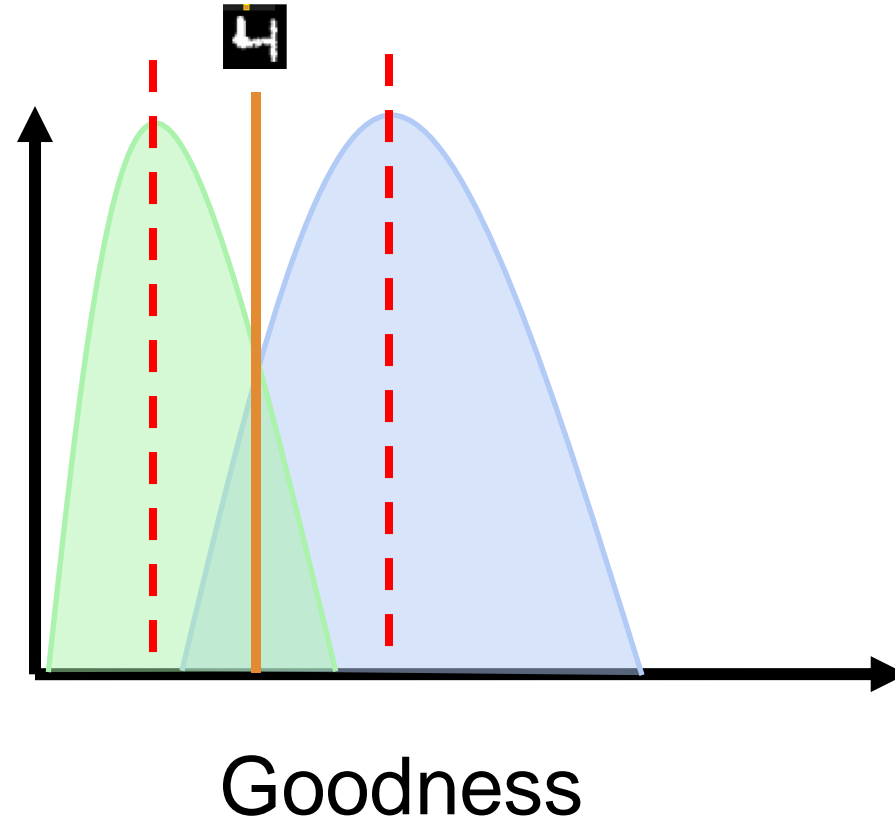
Index: 9



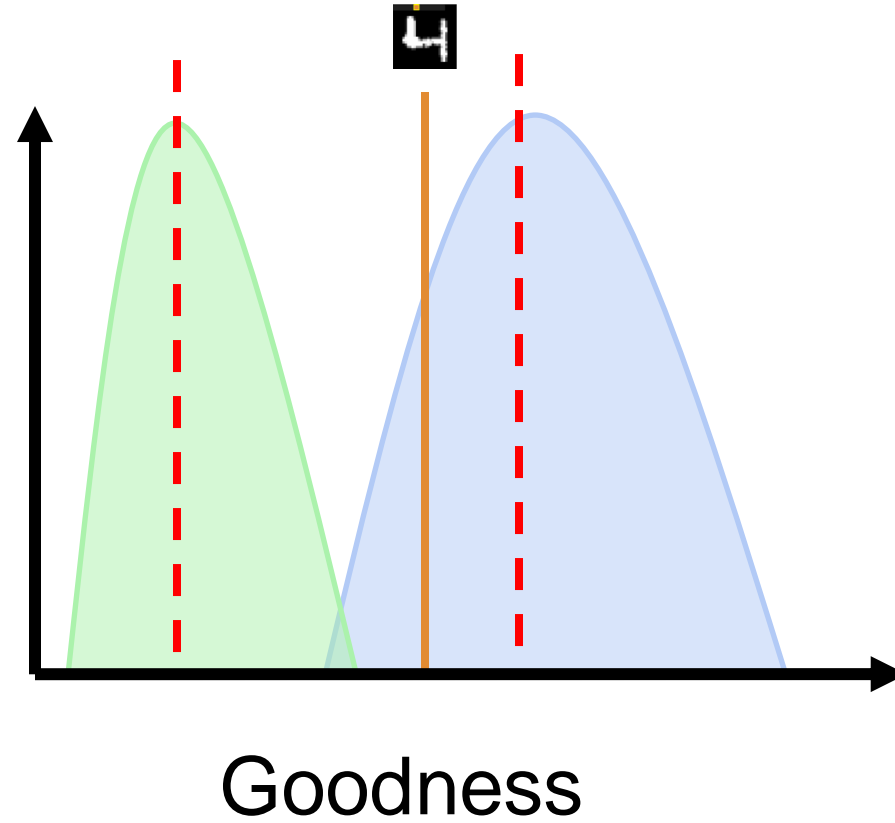
The Goodness is the sum of the square of the activity of each hidden neuron.

Insights From FF

Current Layer

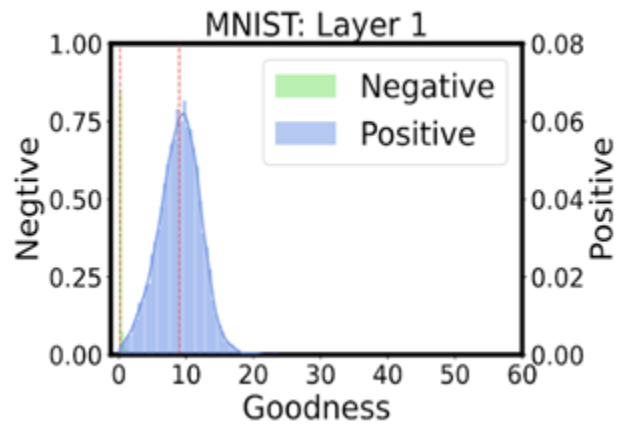


Insights From FF

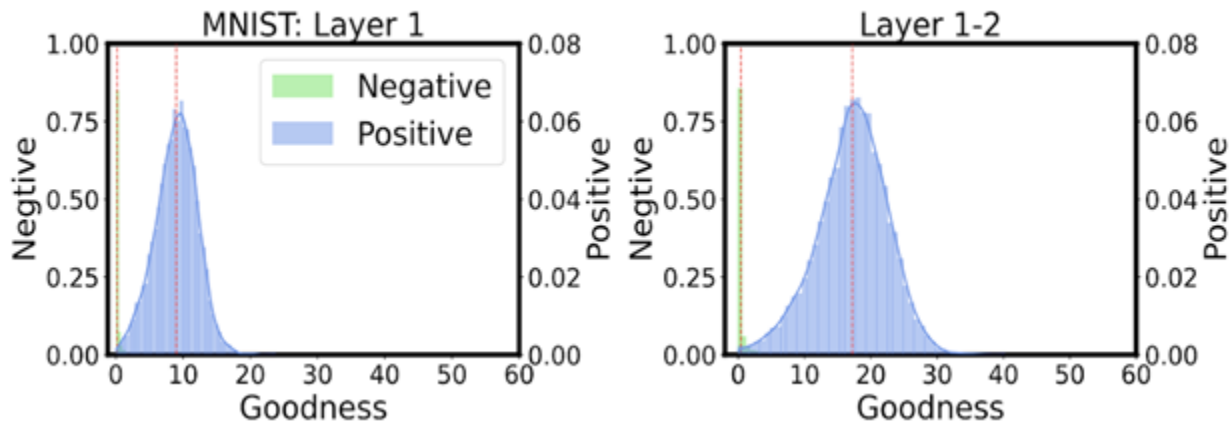


Next Layer

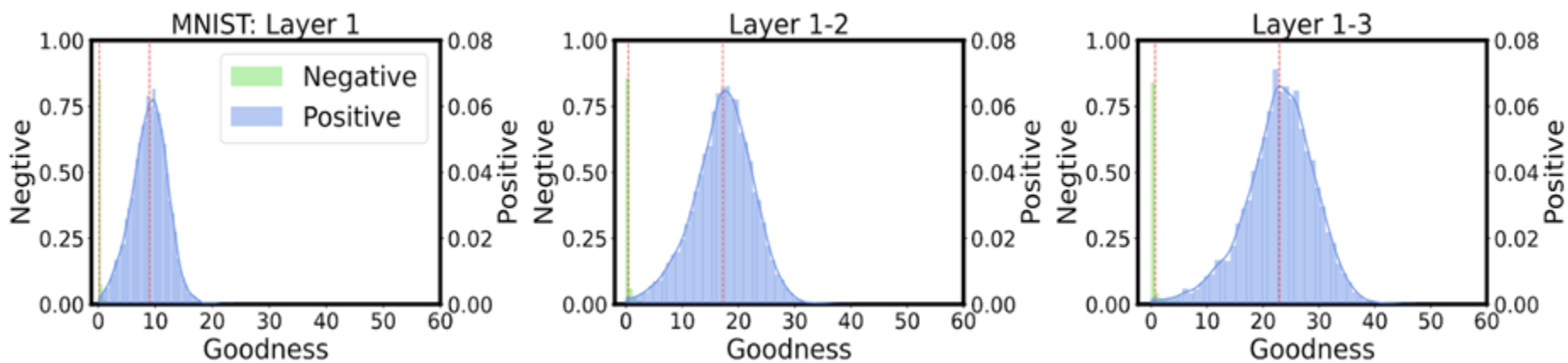
Observation on FF



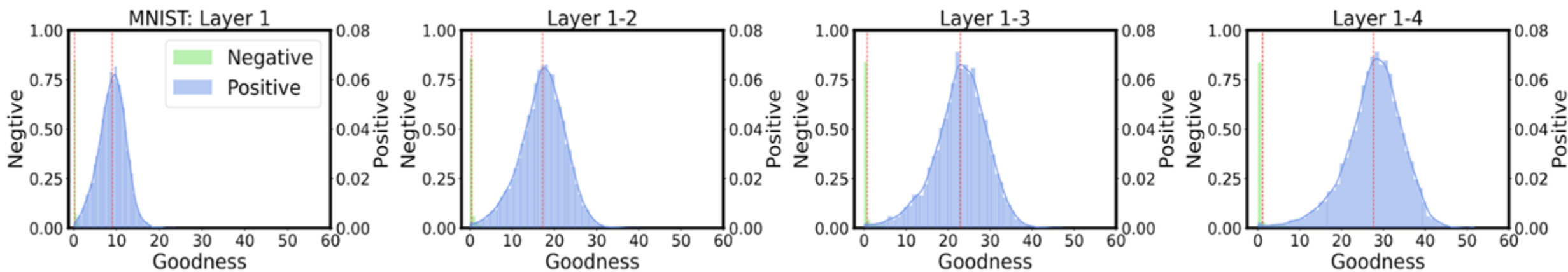
Observation on FF



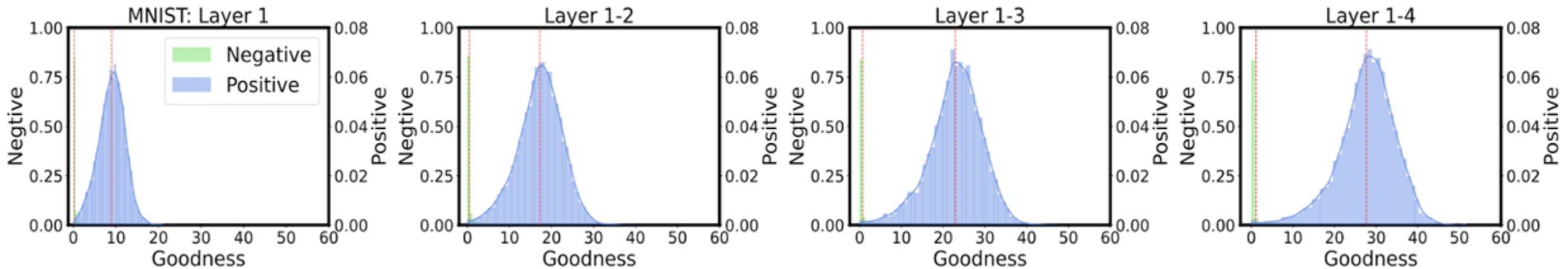
Observation on FF



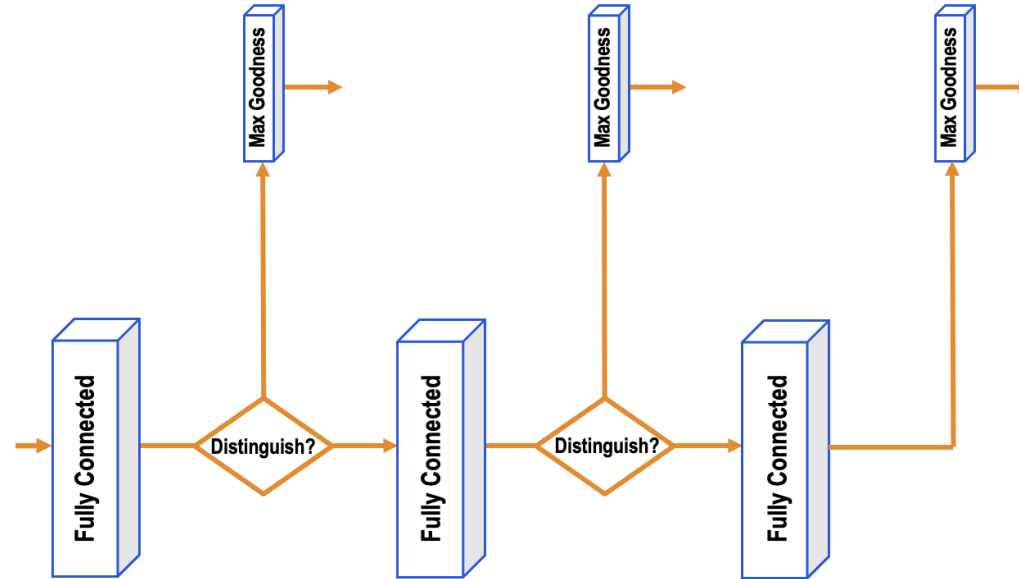
Observation on FF



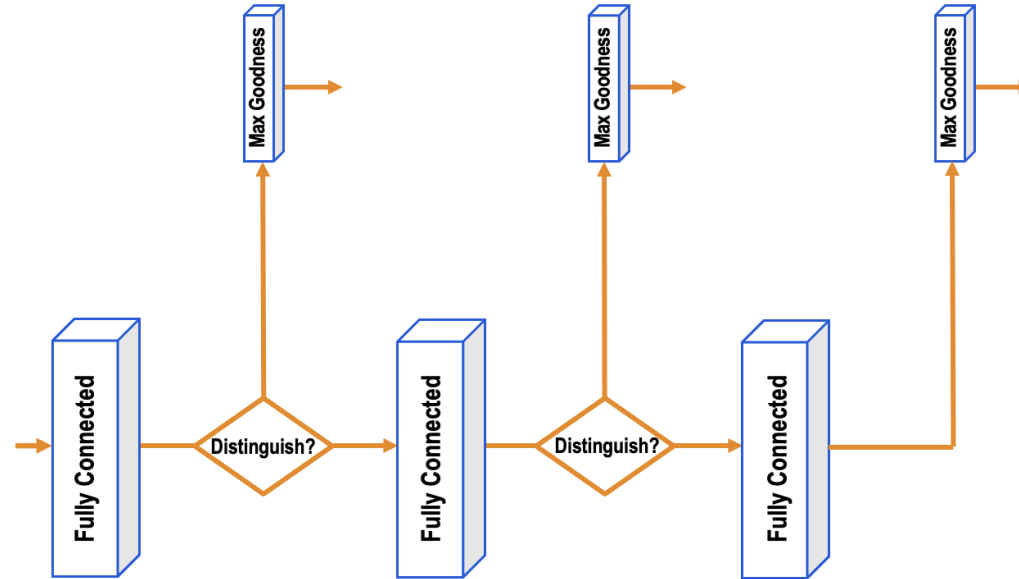
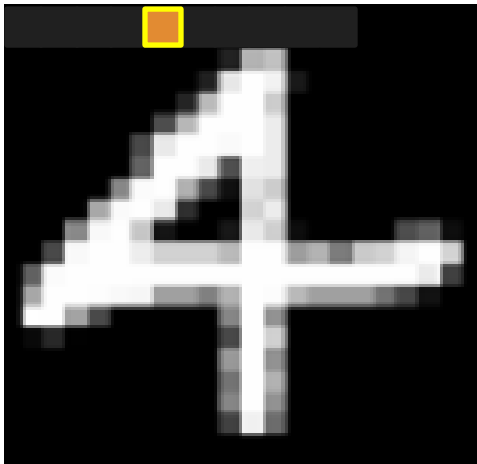
Observation on FF



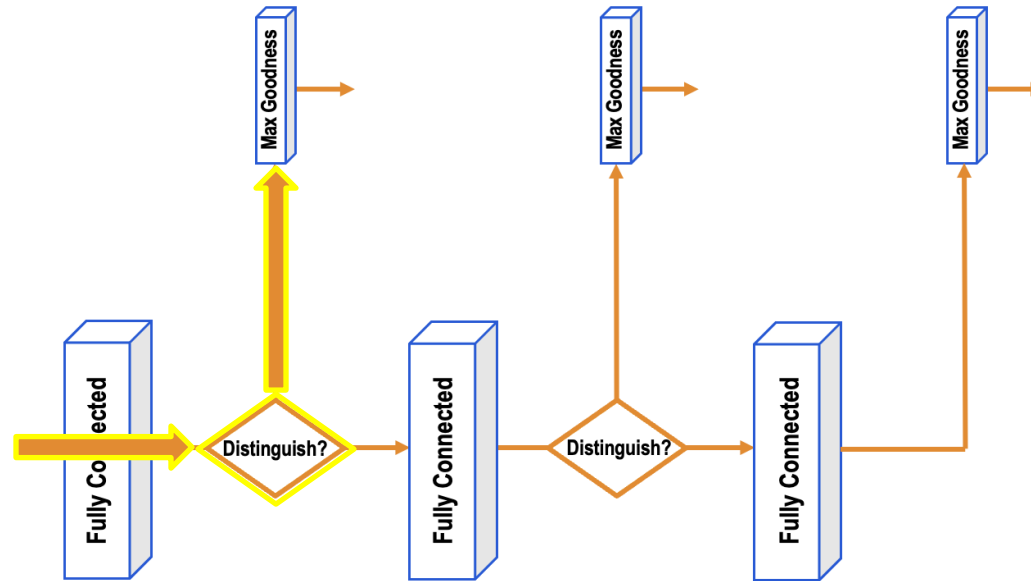
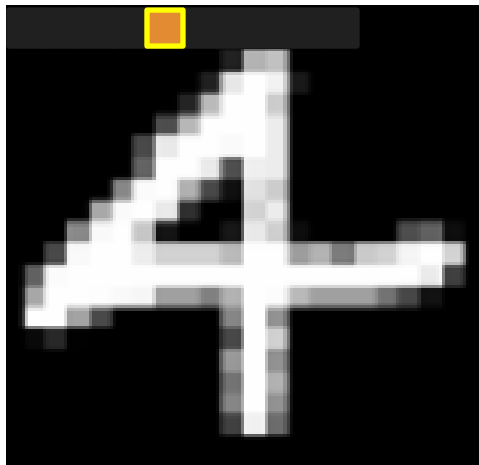
The distance between the mean Goodness of the negative and positive samples/distributions increases as we consider more layers.



Index: 4

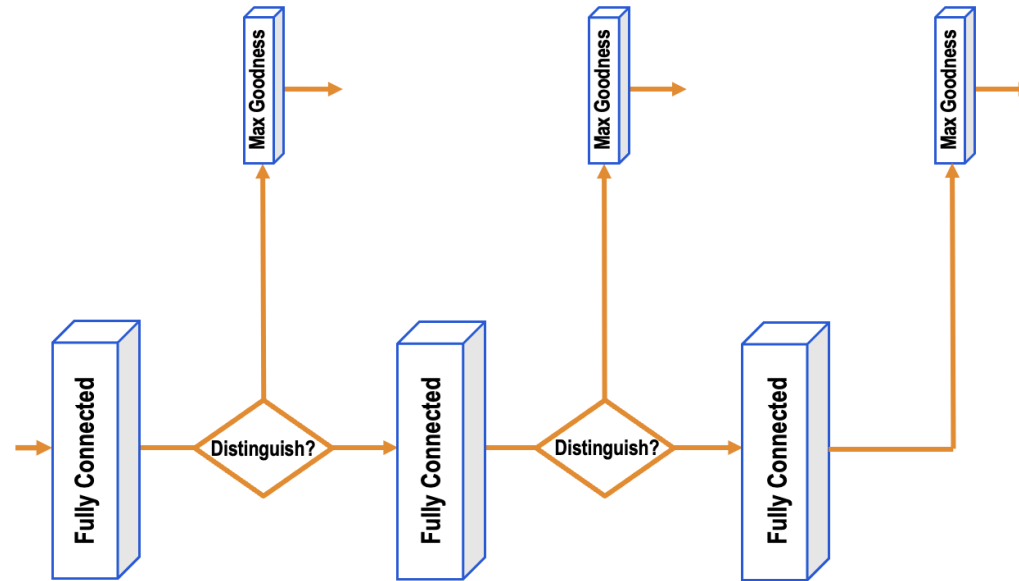
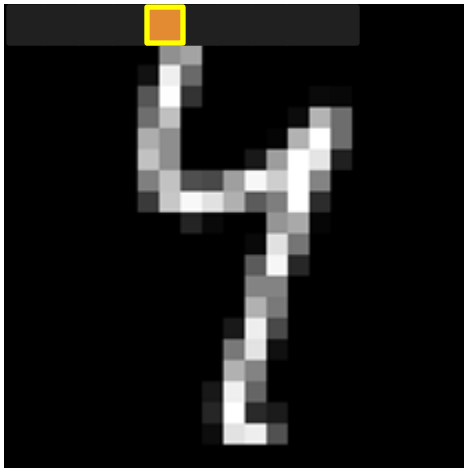


Index: 4

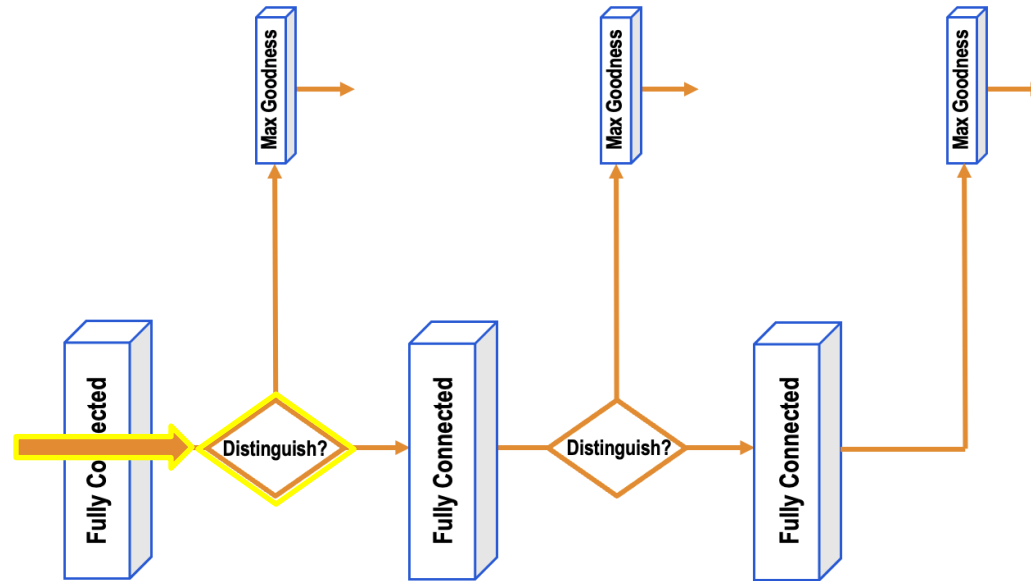
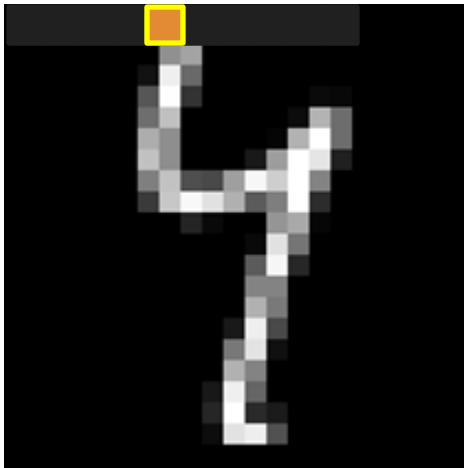


Distinguishable

Index: 4

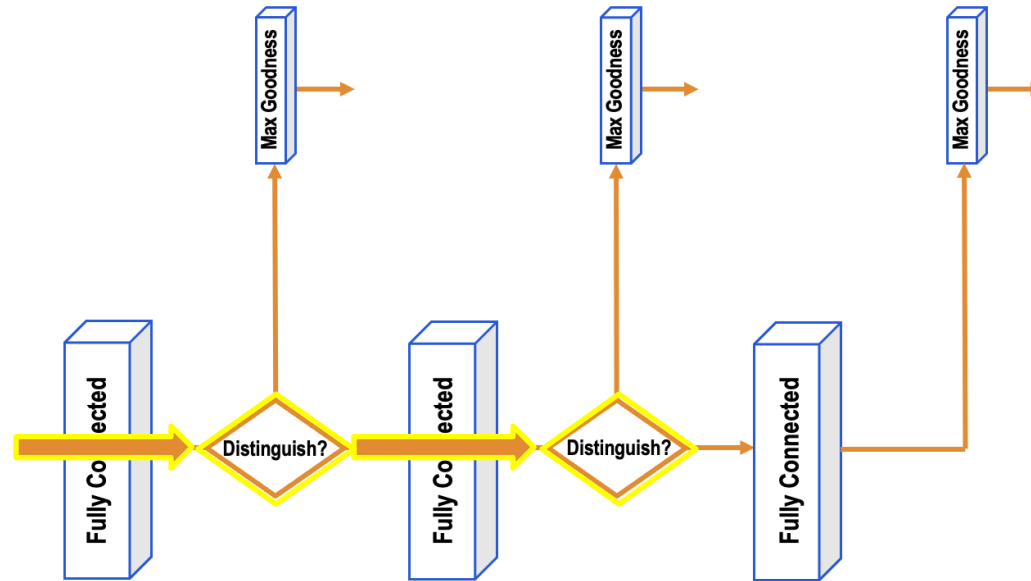
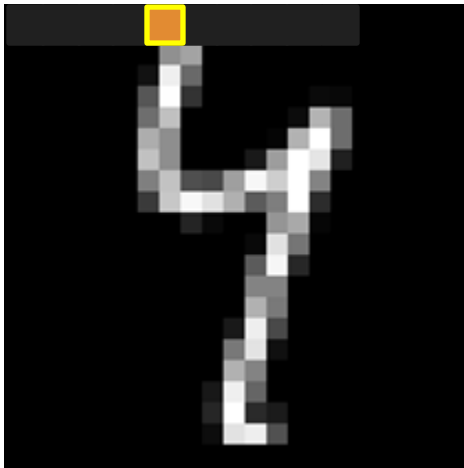


Index: 4

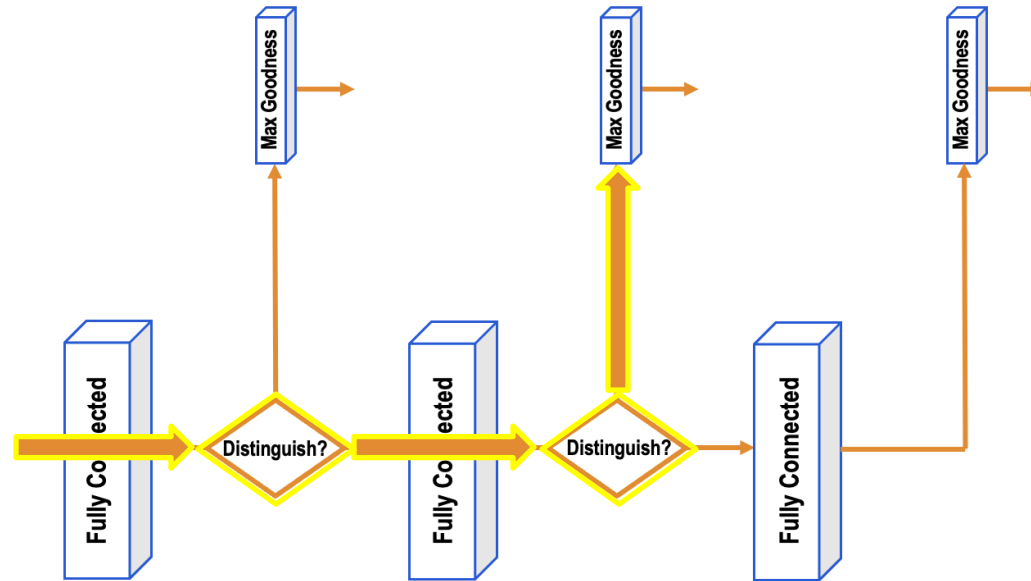
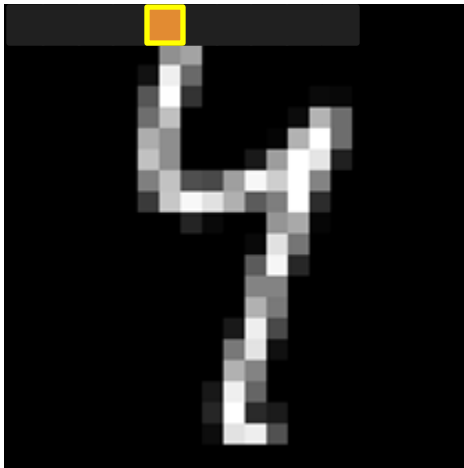


Indistinguishable

Index: 4



Index: 4



Distinguishable



LUND
UNIVERSITY

Evaluation and Results

Dataset and Application



MNIST
Grayscale
Image



CIFAR-10
RGB
Images

Dataset and Application



MNIST
Grayscale
Image



CIFAR-10
RGB
Images



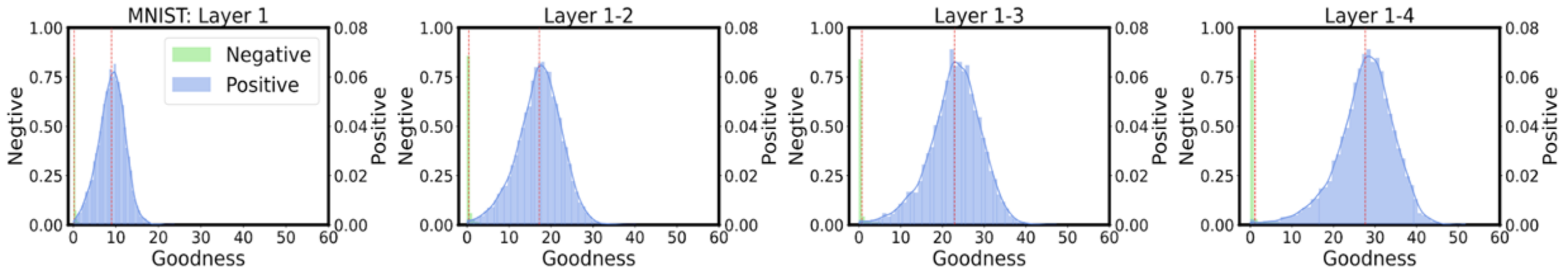
CHB-MIT
Electroencephalogram
(EEG)



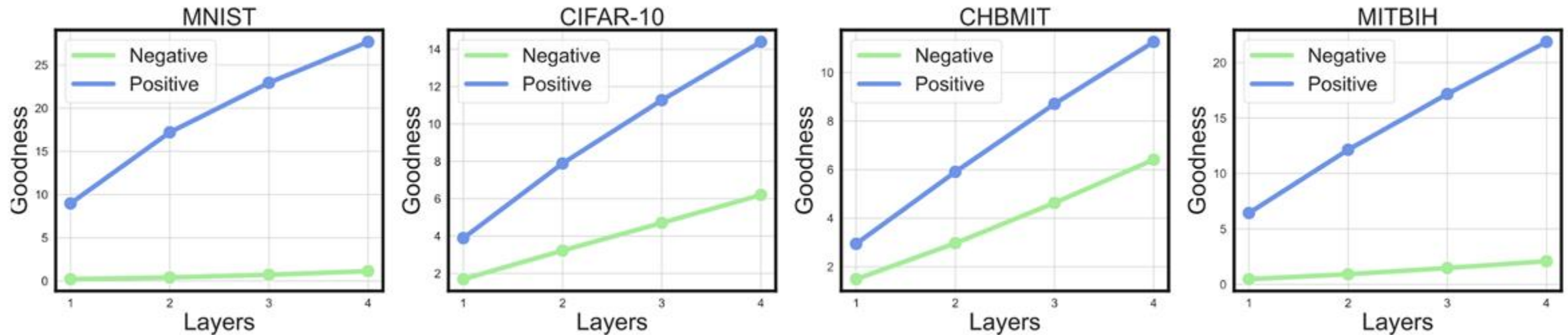
MIT-BIH
Electrocardiogram
(ECG)

**Real-world wearable applications:
Complexity overhead/energy consumption is a major constraint.**

Generalized Insights



Generalized Insights



The distance between the mean Goodness of the negative and positive samples increases as more layers are taken into account.

Classification Performance

Dataset	Error	
	FF	LightFF
MNIST	1.51%	1.51%
CIFAR-10	50.65%	46.05%
CHB-MIT	39.62%	34.83%
MIT-BIH	10.74%	10.07%

LightFF achieves a **comparable** classification error, compared to the Forward-Forward algorithm.

Lightweight Inference for BP

Dataset	Error	
	BP	LightBP
MNIST	1.33%	5.21%
CIFAR-10	43.62%	54.22%
CHB-MIT	25.63%	40.69%
MIT-BIH	8.25%	11.55%

Lightweight Inference for BP

Dataset	Error		
	BP		LightBP
MNIST	1.33%	<	5.21%
CIFAR-10	43.62%	<	54.22%
CHB-MIT	25.63%	<	40.69%
MIT-BIH	8.25%	<	11.55%

Goodness is a good metric for lightweight inference in LightFF.

Inference Complexity

Dataset	MACs	
	FF	LightFF
MNIST	12.94M	2.95M
CIFAR-10	17.30M	10.31M
CHB-MIT	13.39M	8.39M
MIT-BIH	11.90M	1.99M

Inference Complexity

Dataset	MACs			
	FF	LightFF		
MNIST	12.94M	>	2.95M	4.4x
CIFAR-10	17.30M	>	10.31M	1.7x
CHB-MIT	13.39M	>	8.39M	1.6x
MIT-BIH	11.90M	>	1.99M	6.0x

LightFF **improves** computational efficiency, compared to the Forward-Forward algorithm.

Execution Time

Dataset	Execution Time	
	FF	LightFF
MNIST	22.81ms	3.39ms
CIFAR-10	29.57ms	16.38ms
CHB-MIT	4.73ms	2.85ms
MIT-BIH	11.26ms	1.83ms

Execution Time

Dataset	Execution Time			
	FF	LightFF		
MNIST	22.81ms	>	3.39ms	6.7x
CIFAR-10	29.57ms	>	16.38ms	1.8x
CHB-MIT	4.73ms	>	2.85ms	1.7x
MIT-BIH	11.26ms	>	1.83ms	6.2x

LightFF **reduces** the execution time, compared to the Forward-Forward algorithm.

Challenge

Deep Learning
to be **Efficient**



Conclusion

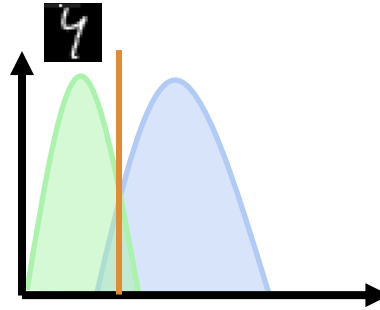
Challenge

Deep Learning
to be **Efficient**



Approach

Lightweight Inference
For Forward-Forward



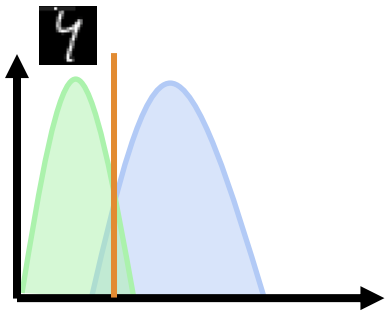
Challenge

Deep Learning
to be **Efficient**



Approach

Lightweight Inference
For Forward-Forward



Performance

Improves Efficiency
Comparable Error

Dataset	MNIST	
	FF	LightFF
Error	1.51%	= 1.51%
Execution Time	22.81ms	> 3.39ms
MACs	12.94M	> 2.95M

Conclusion

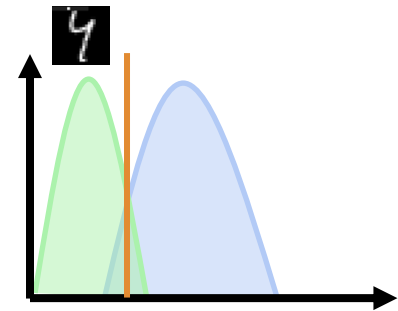
Challenge

Deep Learning
to be **Efficient**



Approach

Lightweight Inference
For Forward-Forward



Performance

Improves Efficiency
Comparable Error

Dataset	MNIST	
	FF	LightFF
Error	1.51%	= 1.51%
Execution Time	22.81ms	> 3.39ms
MACs	12.94M	> 2.95M



Paper

Wednesday Outreach Activities

Thank you!



Code