

Modern Machine Learning

The Good and **the Bad**

Amartya Sanyal

Postdoc @ Max-Planck Institute for Intelligent Systems, Tübingen

(Incoming) Tenure Track Assistant Professor in Machine Learning, Department of Computer Science

(Affiliated) Assistant Professor, Department of Mathematics,
University of Copenhagen

CISPA

MAX PLANCK INSTITUTE
FOR INTELLIGENT SYSTEMS



ETH zürich



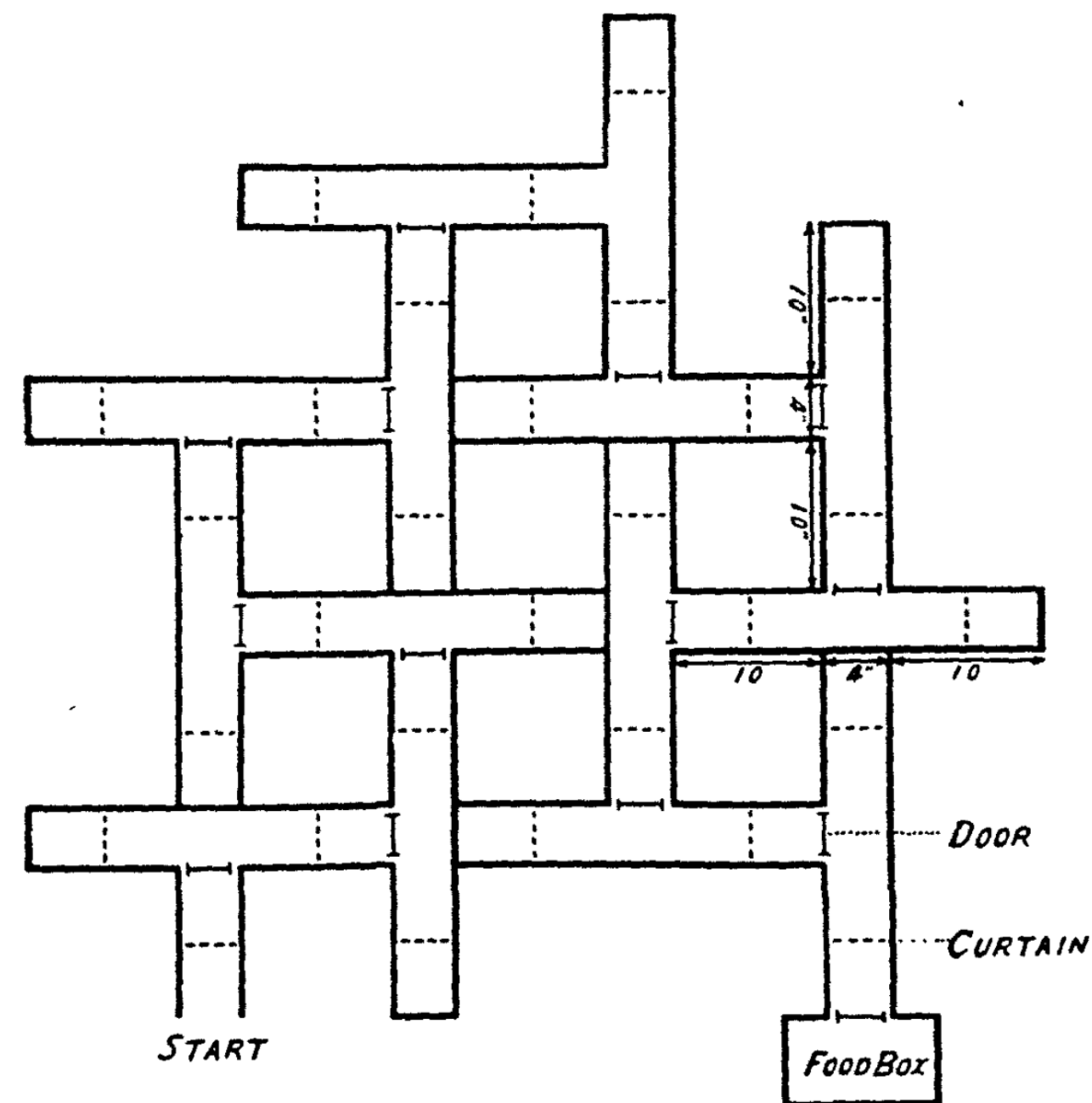
Learning

Learning

“Learning is any process by which a system improves performance from experience.” - Herbert Simon

Learning

“Learning is any process by which a system improves performance from experience.” - Herbert Simon



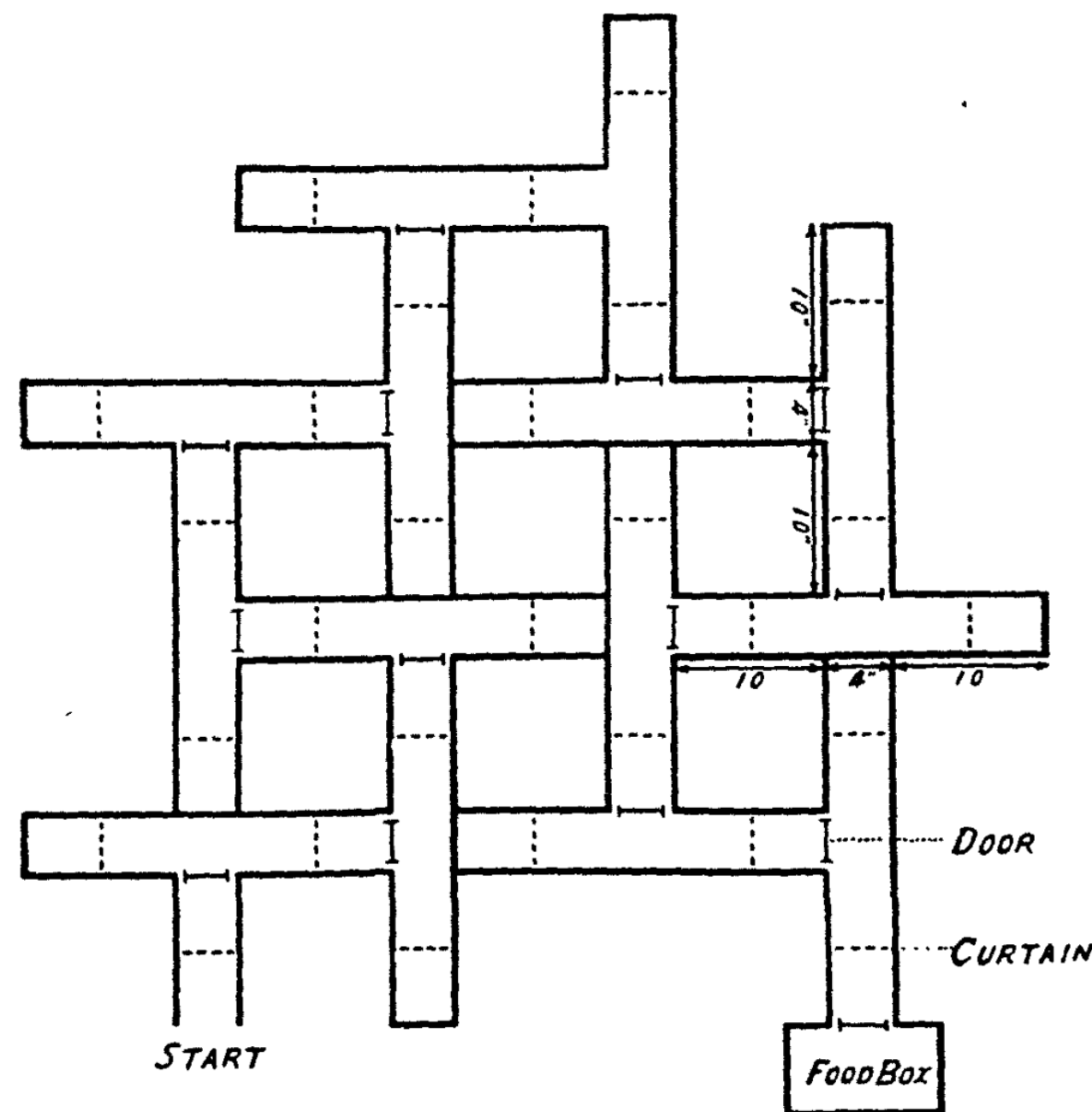
Plan of maze
14-Unit T-Alley Maze

FIG. 1

(From M. H. Elliott, The effect of change of reward on the maze performance of rats. *Univ. Calif. Publ. Psychol.*, 1928, 4, p. 20.)

Learning

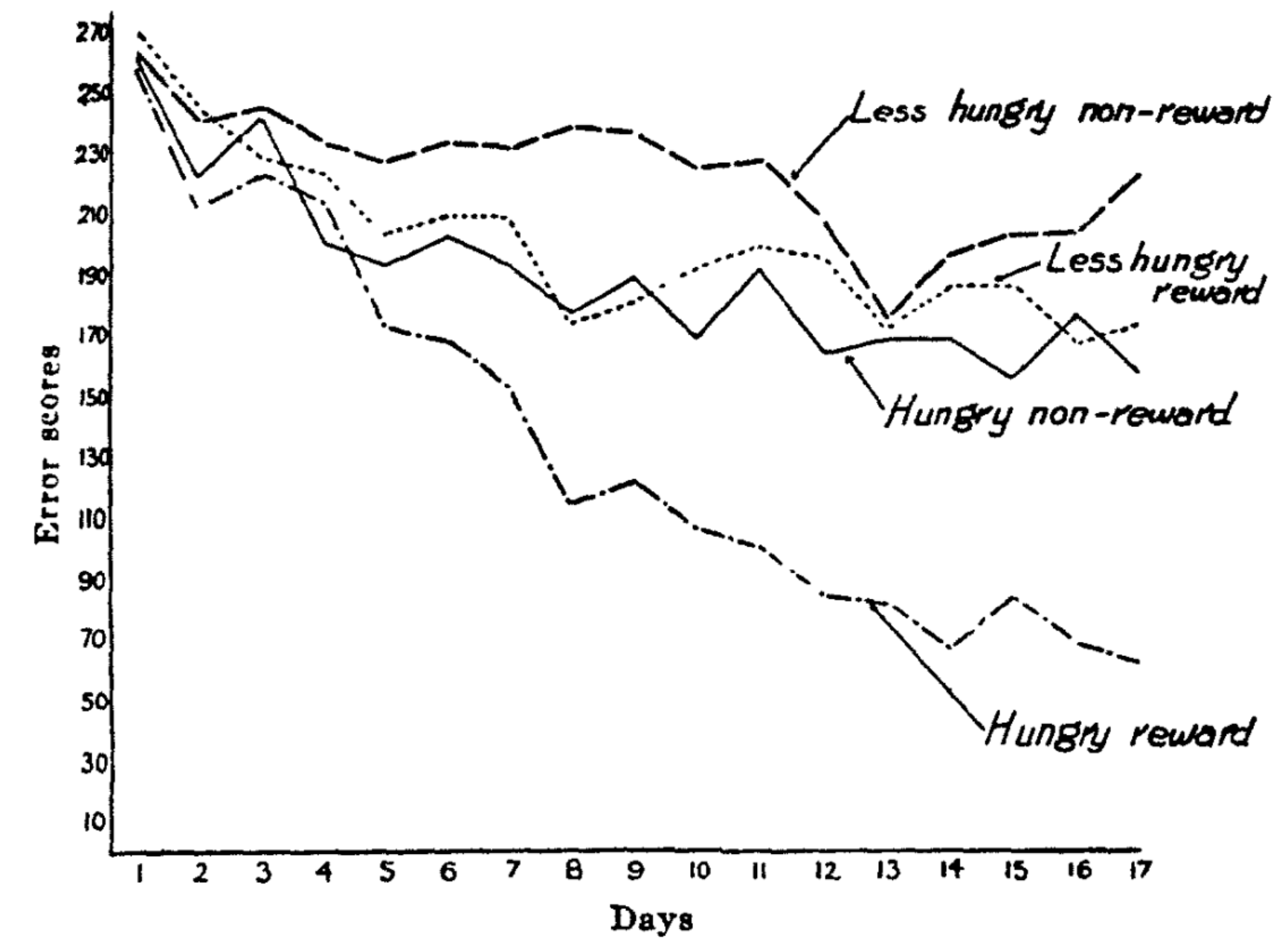
“Learning is any process by which a system improves performance from experience.” - Herbert Simon



Plan of maze
14-Unit T-Alley Maze

FIG. 1

(From M. H. Elliott, The effect of change of reward on the maze performance of rats. *Univ. Calif. Publ. Psychol.*, 1928, 4, p. 20.)



Error curves for four groups, 36 rats.

FIG. 3

(From E. C. Tolman and C. H. Honzik, Degrees of hunger, reward and non-reward, and maze learning in rats. *Univ. Calif. Publ. Psychol.*, 1930, 4, No. 16, p. 246. A maze identical with the alley maze shown in Fig. 1 was used.)

Machine Learning (ML)

“Learning is any process by which a system improves performance from experience.” - Herbert Simon

Machine Learning (ML)

“Learning is any process by which a system improves performance from experience.” - Herbert Simon

“Machine Learning is the study of algorithms that improve their performance from past experience” - Tom Mitchell

Example: Thermostat

Example: Thermostat



Traditional Decision Making

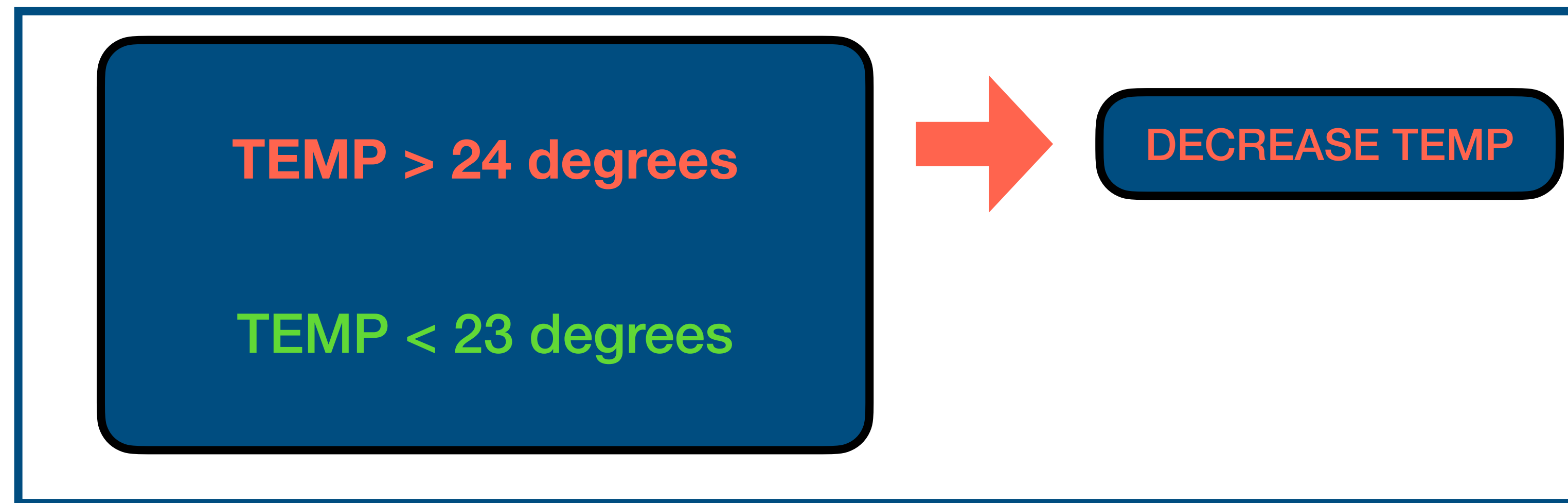
Example: Thermostat

TEMP > 24 degrees

TEMP < 23 degrees

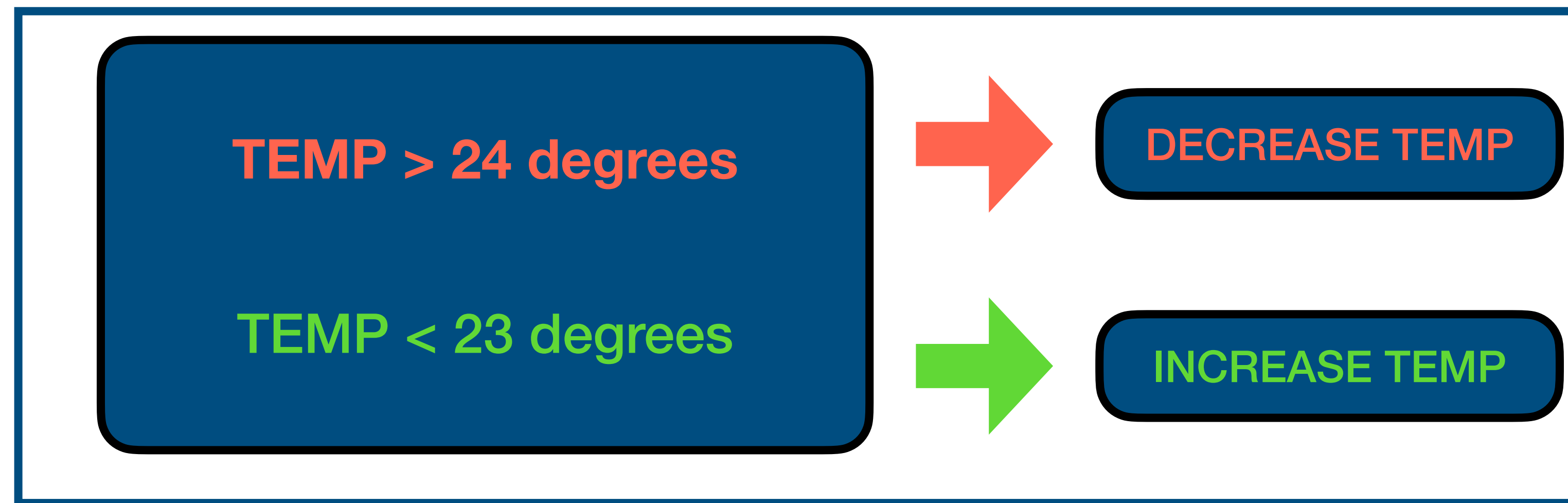
Traditional Decision Making

Example: Thermostat



Traditional Decision Making

Example: Thermostat



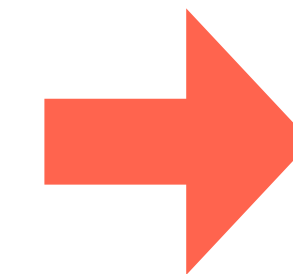
Traditional Decision Making

Example: ML for Thermostat

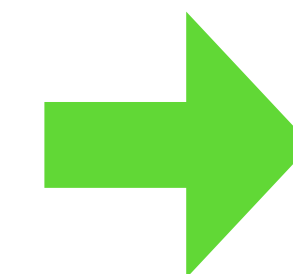
Machine Learning
Algorithm

TEMP > 24 degrees

TEMP < 23 degrees



DECREASE TEMP

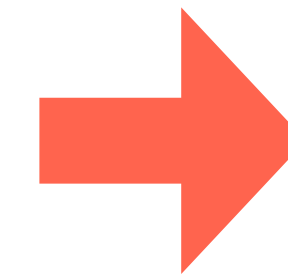


INCREASE TEMP

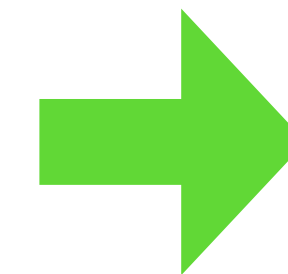
Traditional Decision Making

Example: ML for Thermostat

Machine Learning
Algorithm

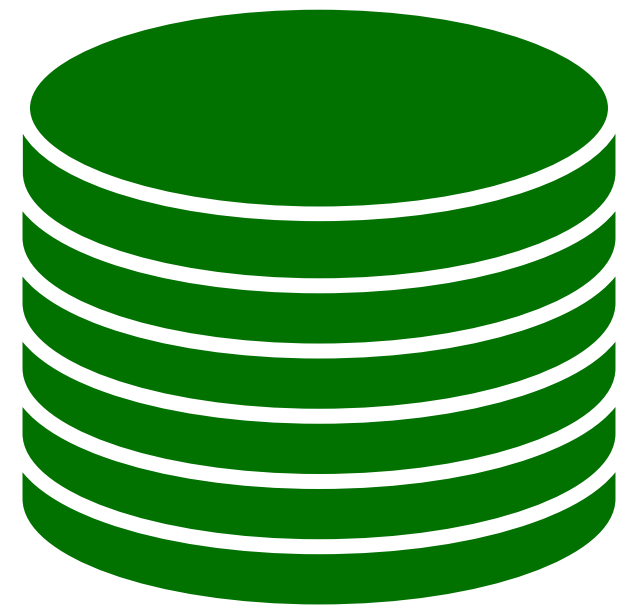


DECREASE TEMP



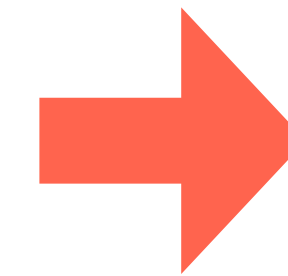
INCREASE TEMP

Example: ML for Thermostat

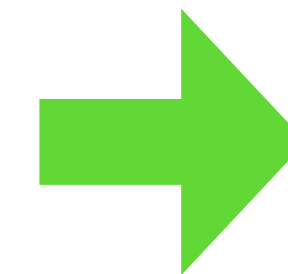


Past Data on

Machine Learning
Algorithm

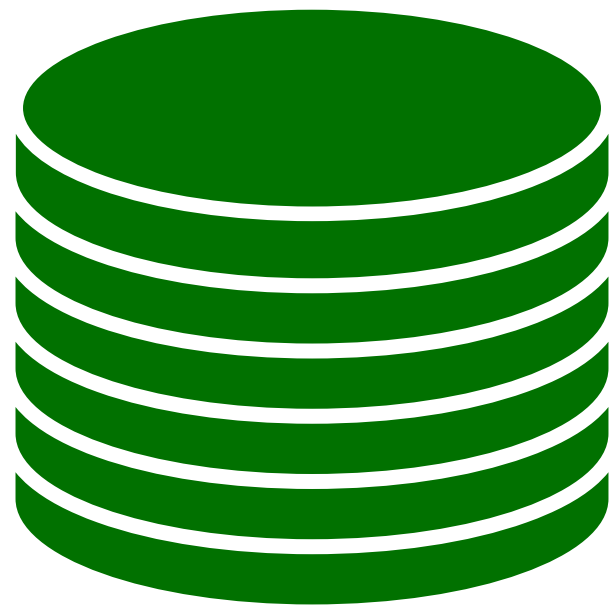


DECREASE TEMP



INCREASE TEMP

Example: ML for Thermostat

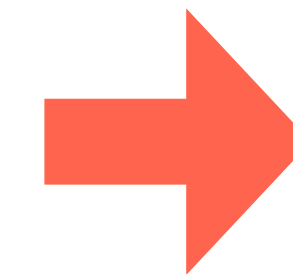


Past Data on

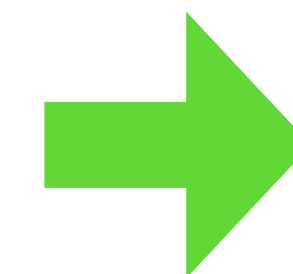
- Past Manual Temp Change



Machine Learning
Algorithm

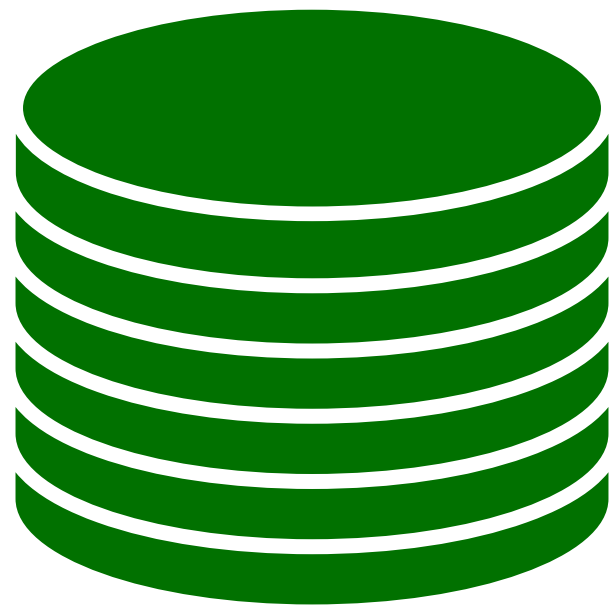


DECREASE TEMP



INCREASE TEMP

Example: ML for Thermostat

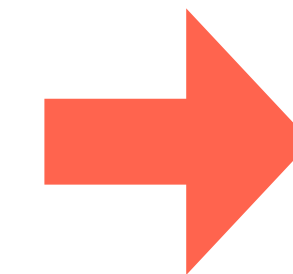


Past Data on

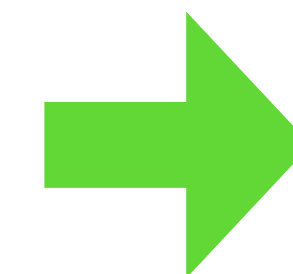
- Past Manual Temp Change
- Time of Day



Machine Learning
Algorithm

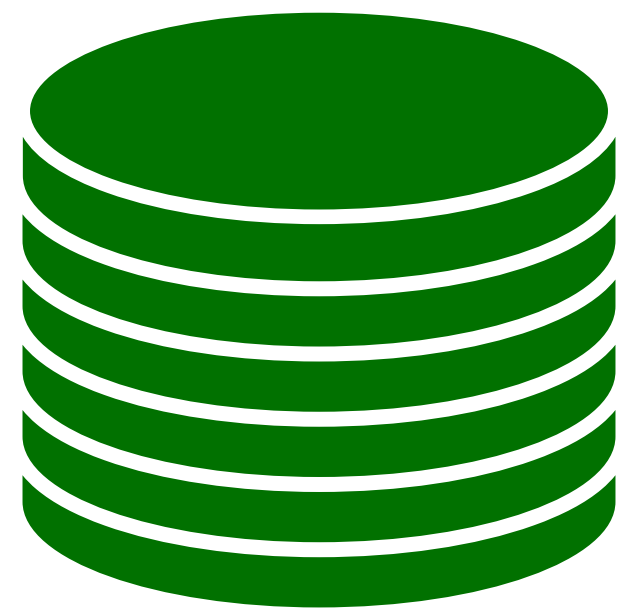


DECREASE TEMP



INCREASE TEMP

Example: ML for Thermostat

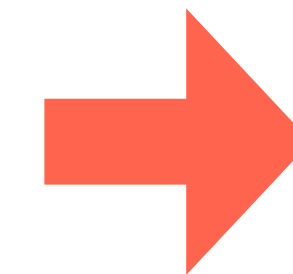


Past Data on

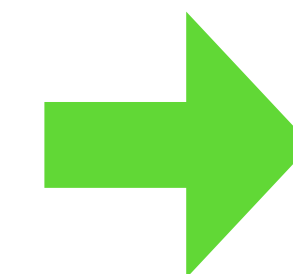
- Past Manual Temp Change
- Time of Day
- Humidity



Machine Learning
Algorithm

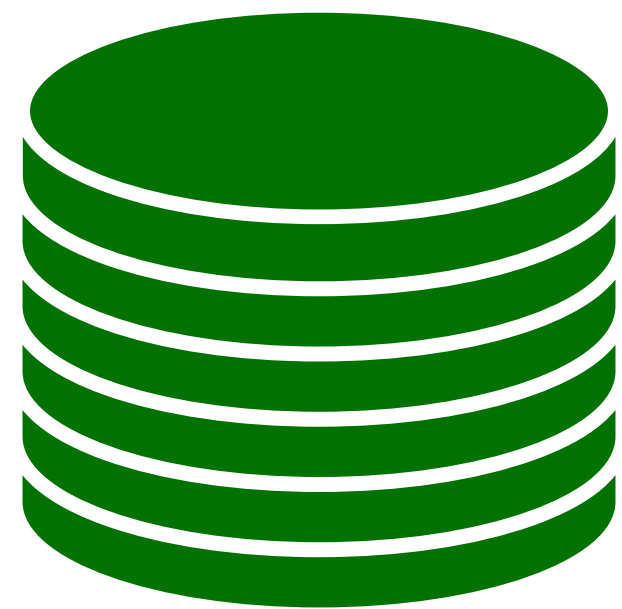


DECREASE TEMP



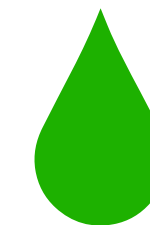
INCREASE TEMP

Example: ML for Thermostat

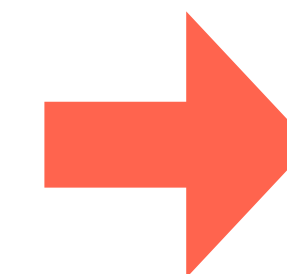


Past Data on

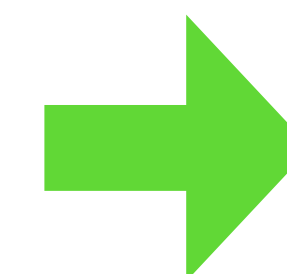
- Past Manual Temp Change
- Time of Day
- Humidity
- Room Temperature



Machine Learning
Algorithm

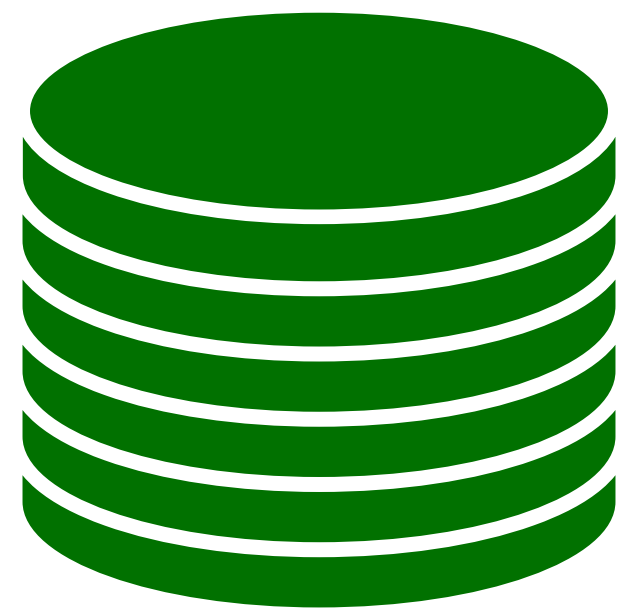


DECREASE TEMP



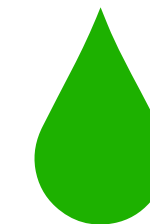
INCREASE TEMP

Example: ML for Thermostat

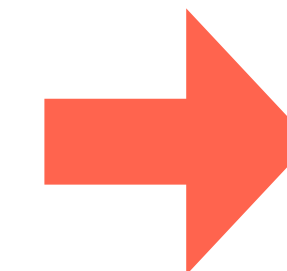


Past Data on

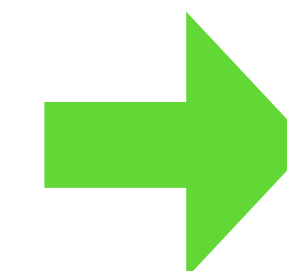
- Past Manual Temp Change
- Time of Day
- Humidity
- Room Temperature
- Number of people in room



Machine Learning
Algorithm

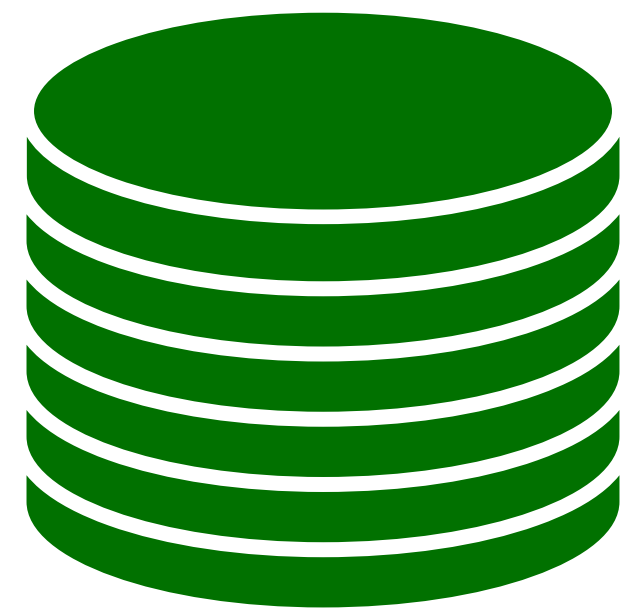


DECREASE TEMP



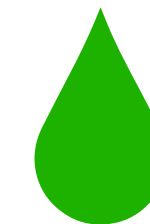
INCREASE TEMP

Example: ML for Thermostat

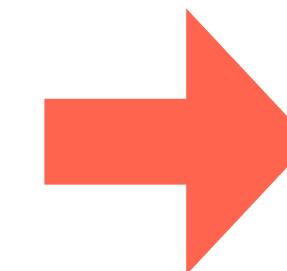


Past Data on

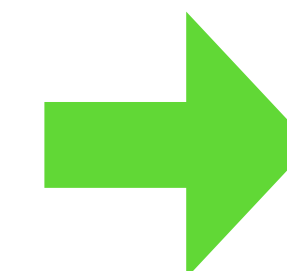
- Past Manual Temp Change
- Time of Day
- Humidity
- Room Temperature
- Number of people in room



Machine Learning
Algorithm

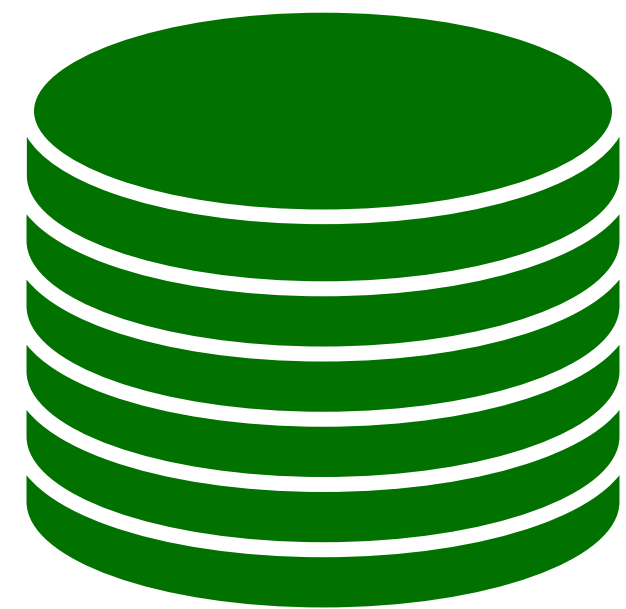


DECREASE TEMP



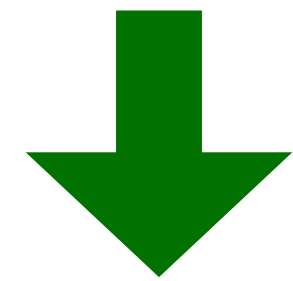
INCREASE TEMP

Example: ML for Thermostat

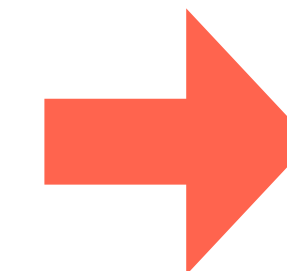
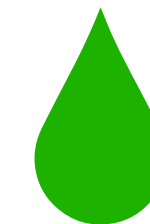


Past Data on

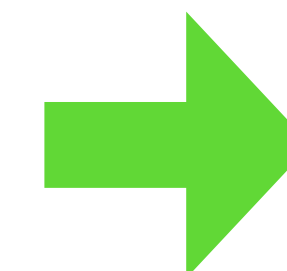
- Past Manual Temp Change
- Time of Day
- Humidity
- Room Temperature
- Number of people in room



Machine Learning
Algorithm

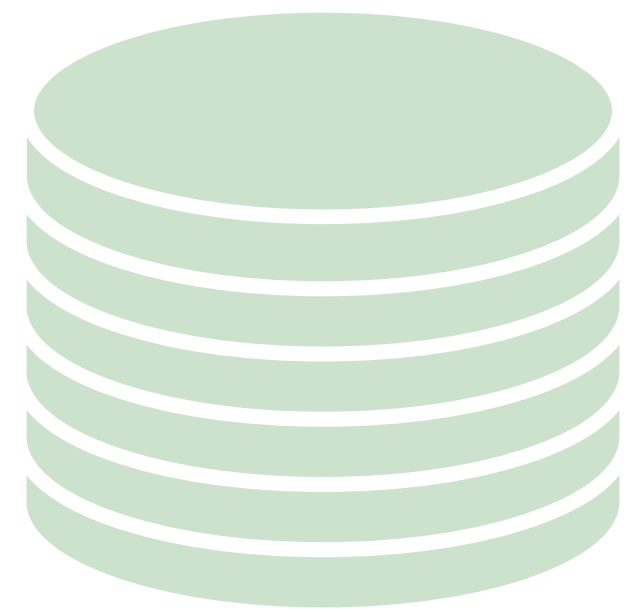


DECREASE TEMP



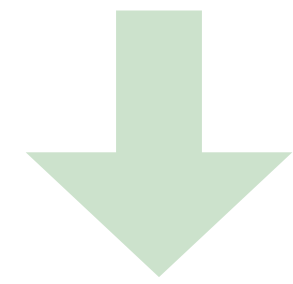
INCREASE TEMP

Example: ML for Thermostat

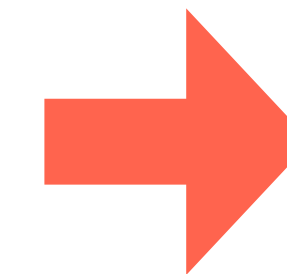
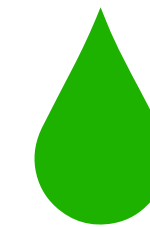


Past Data on

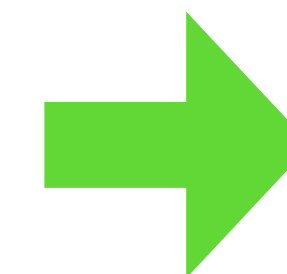
- Past Manual Temp Change
- Time of Day
- Humidity
- Room Temperature
- Number of people in room



**Machine Learning
Algorithm**

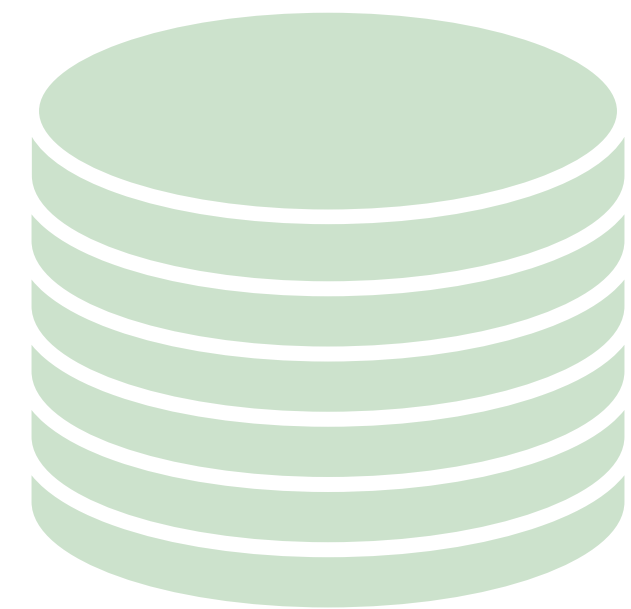


DECREASE TEMP



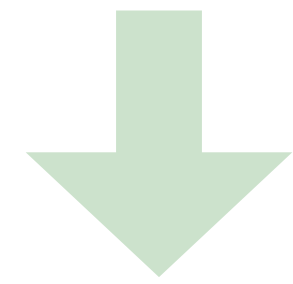
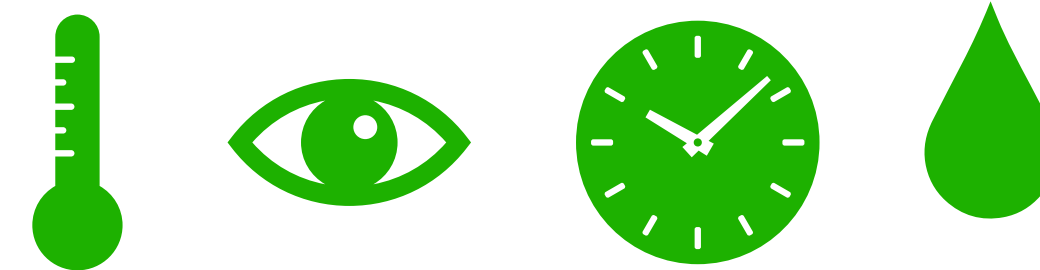
INCREASE TEMP

Example: ML for Thermostat

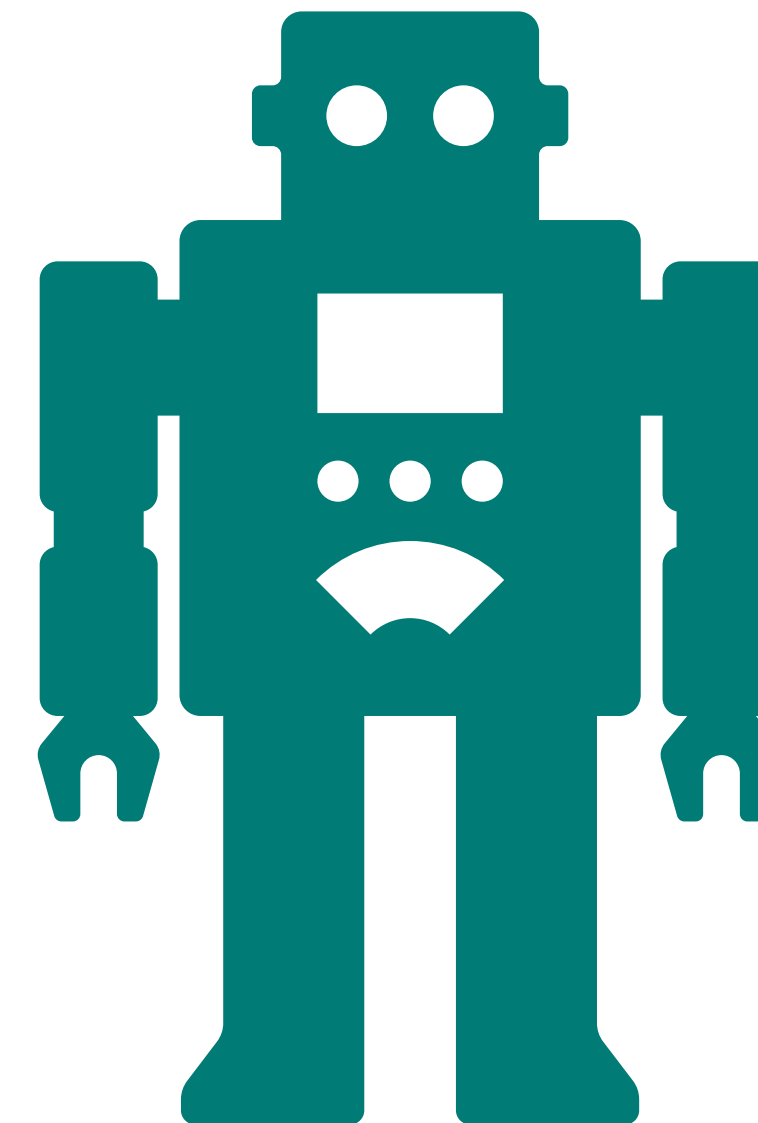
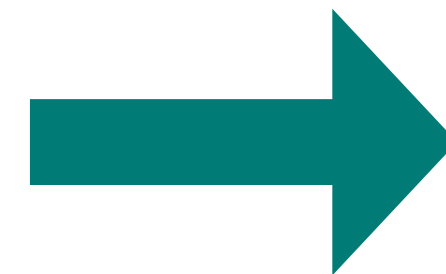


Past Data on

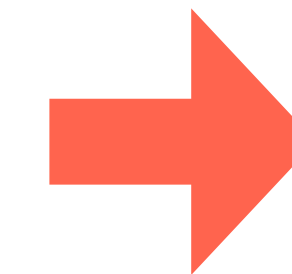
- Past Manual Temp Change
- Time of Day
- Humidity
- Room Temperature
- Number of people in room



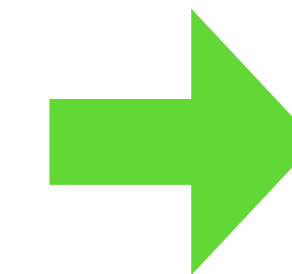
**Machine Learning
Algorithm**



ML Model

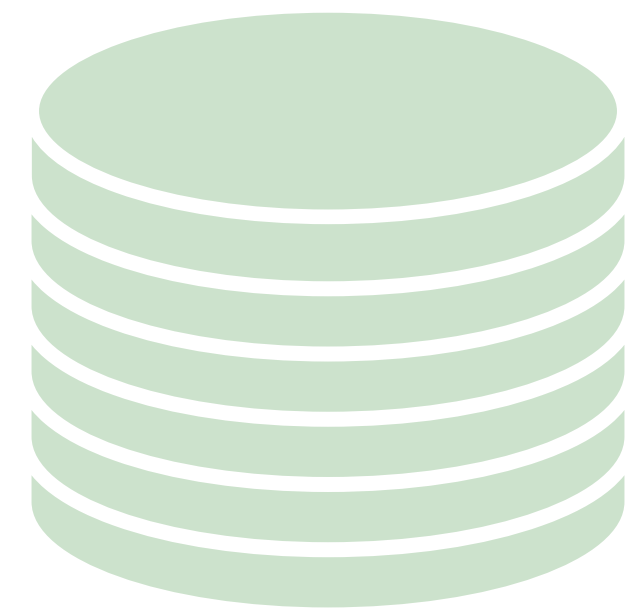


DECREASE TEMP



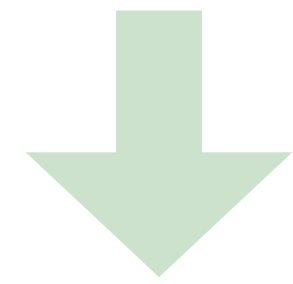
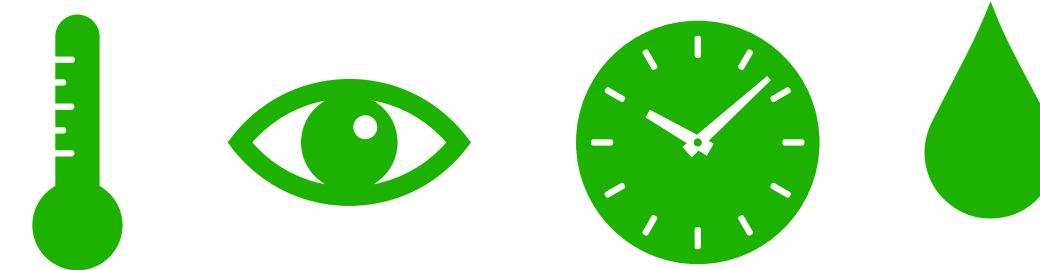
INCREASE TEMP

Example: ML for Thermostat

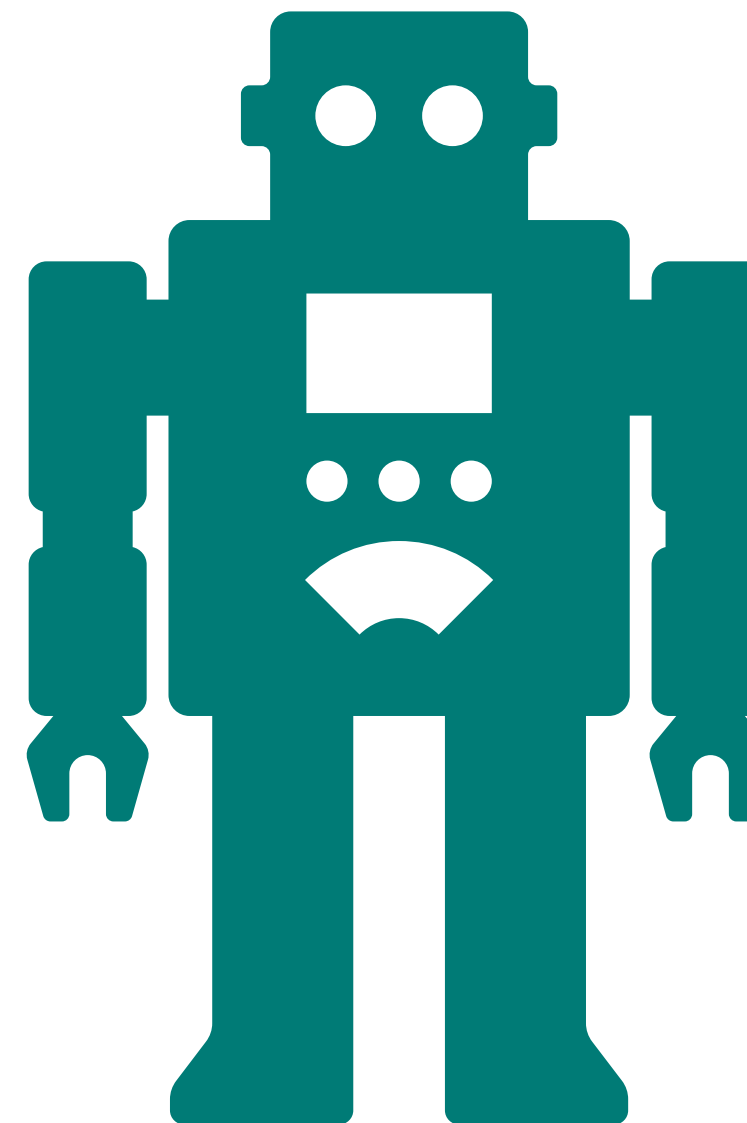
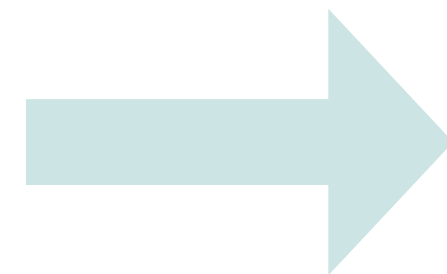


Past Data on

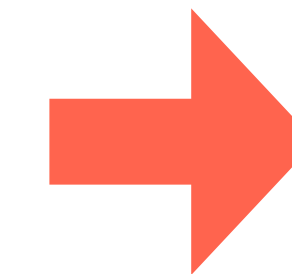
- Past Manual Temp Change
- Time of Day
- Humidity
- Room Temperature
- Number of people in room



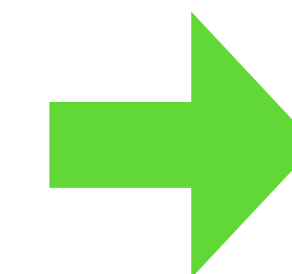
Machine Learning
Algorithm



ML Model

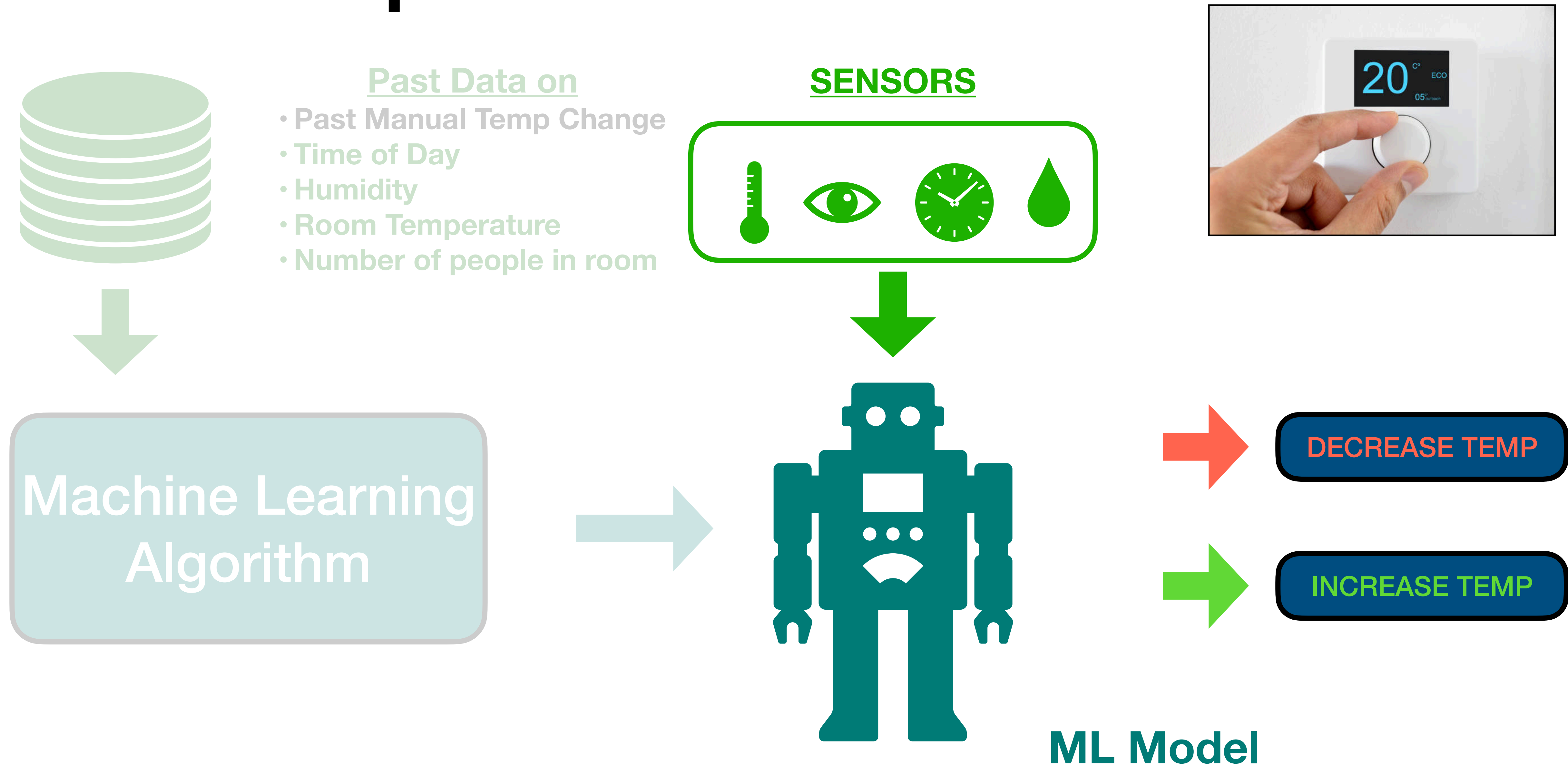


DECREASE TEMP



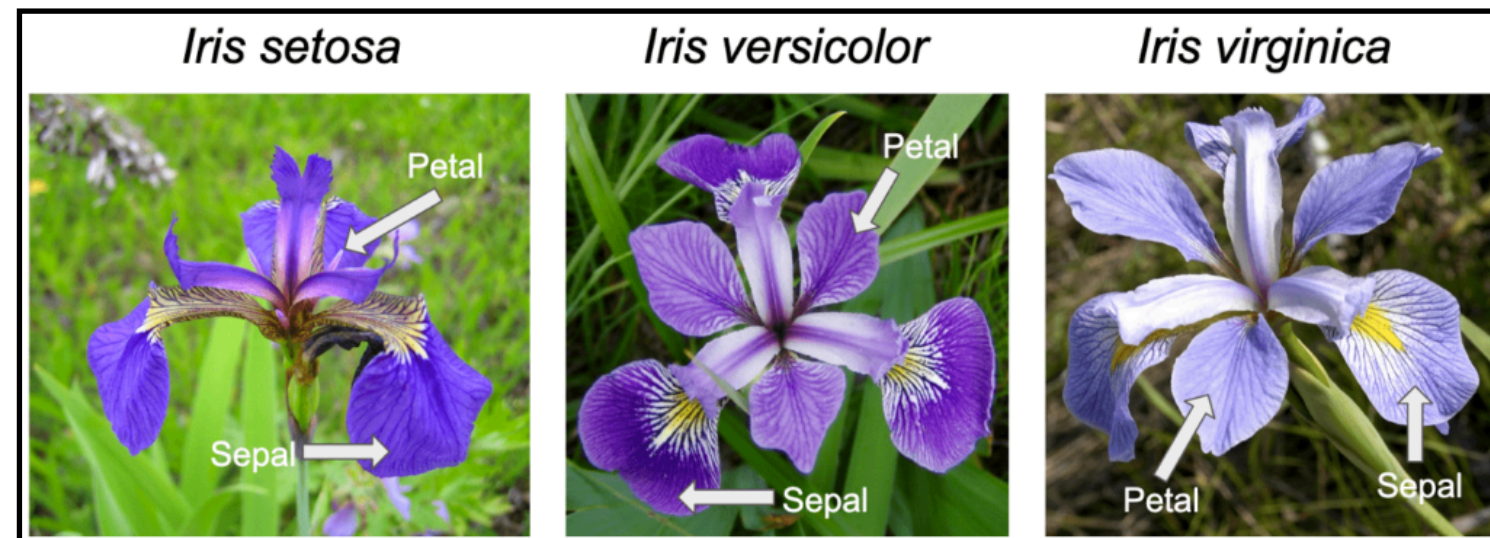
INCREASE TEMP

Example: ML for Thermostat

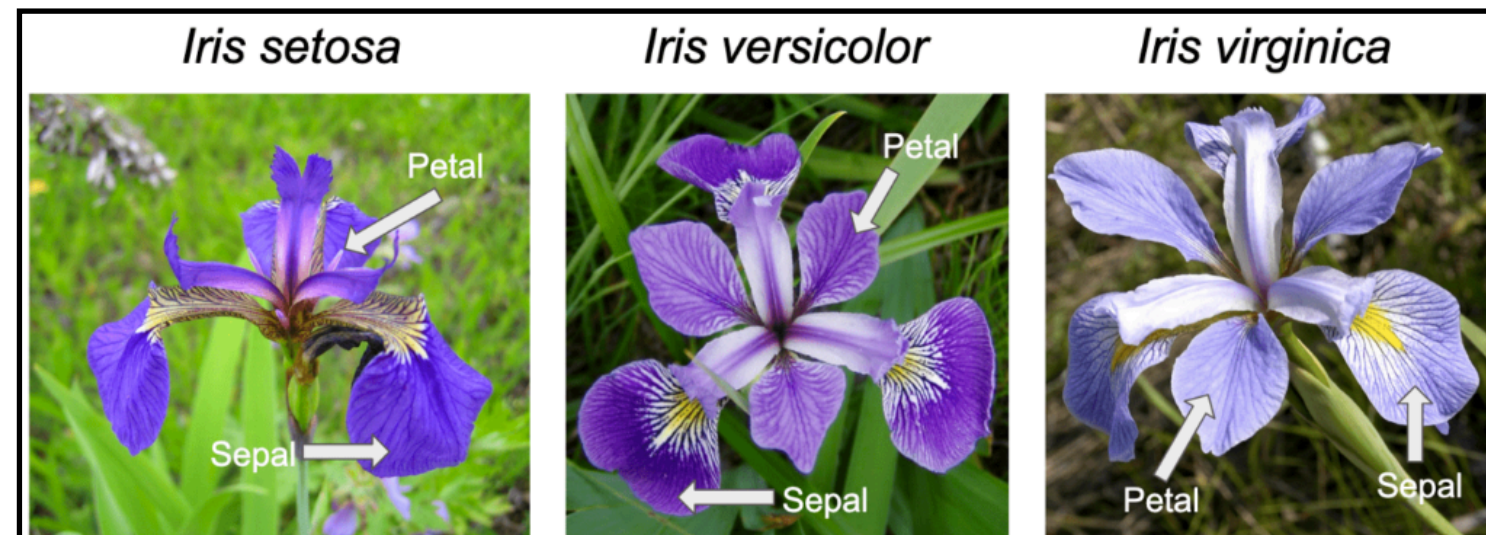


Application: Multi-Class Classification

Application: Multi-Class Classification



Application: Multi-Class Classification



Samples
(instances, observations)

	Sepal length	Sepal width	Petal length	Petal width	Class label
1	5.1	3.5	1.4	0.2	Setosa
2	4.9	3.0	1.4	0.2	Setosa
...					
50	6.4	3.5	4.5	1.2	Versicolor
...					
150	5.9	3.0	5.0	1.8	Virginica

Features
(attributes, measurements, dimensions)

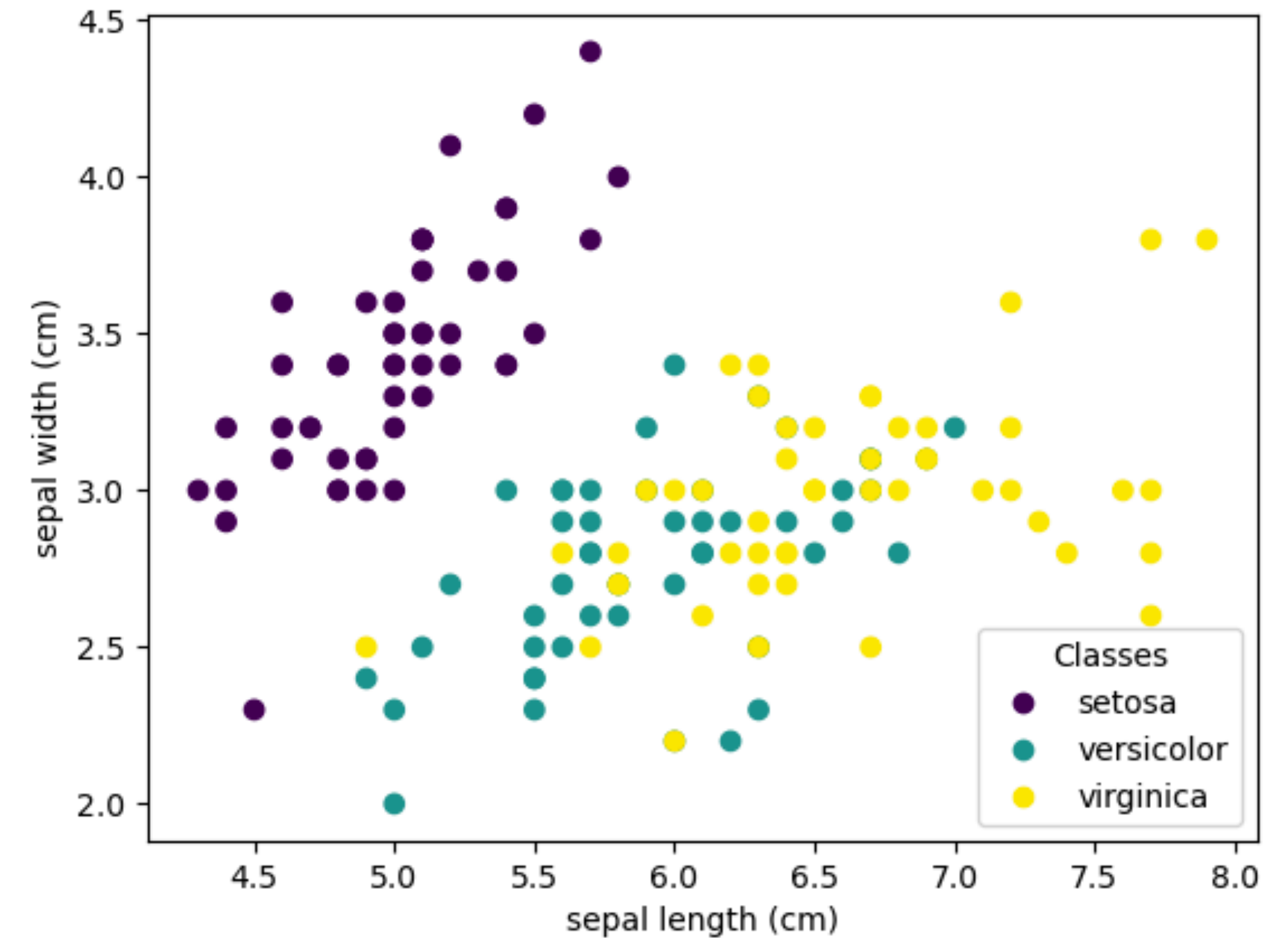
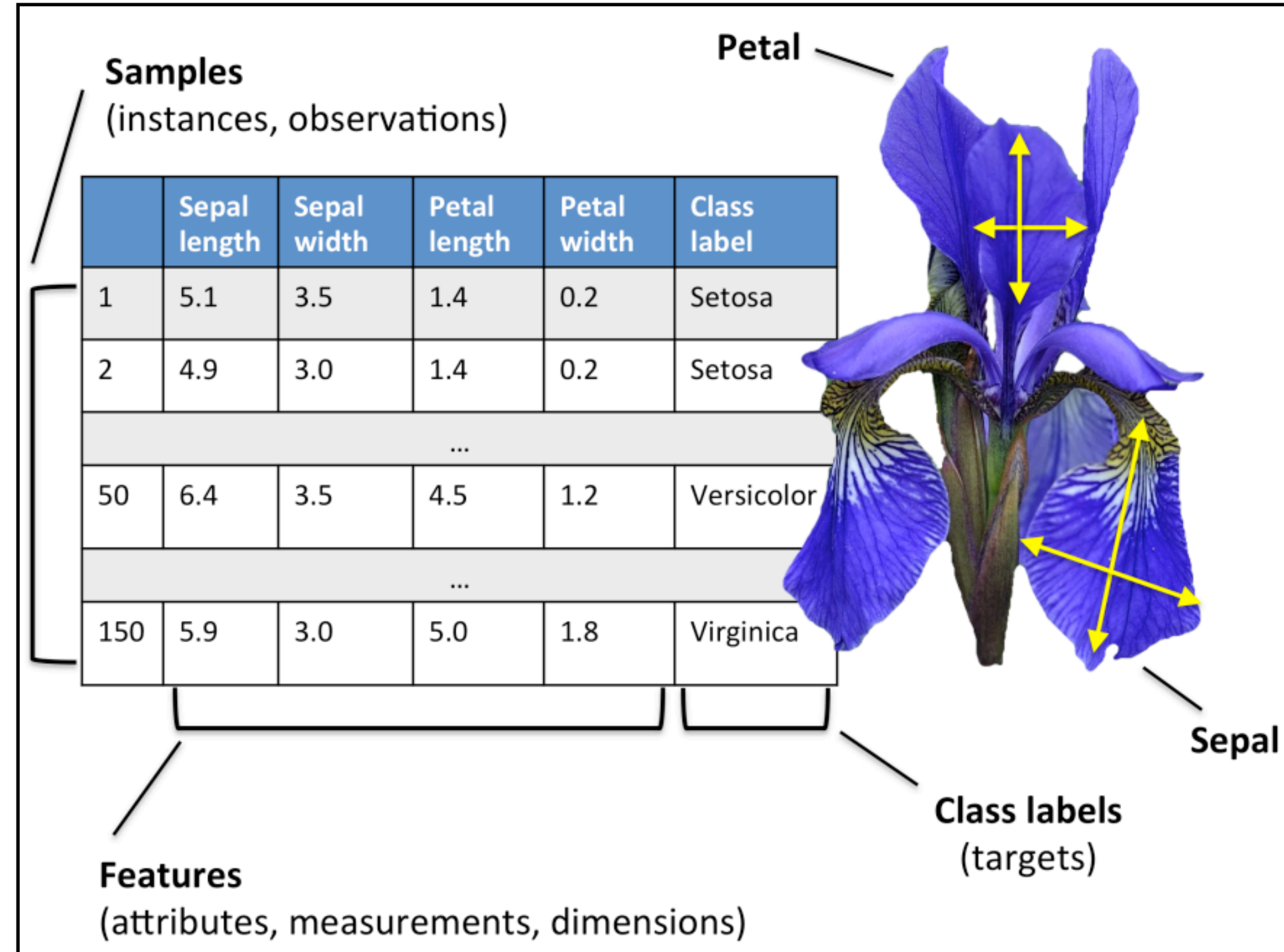
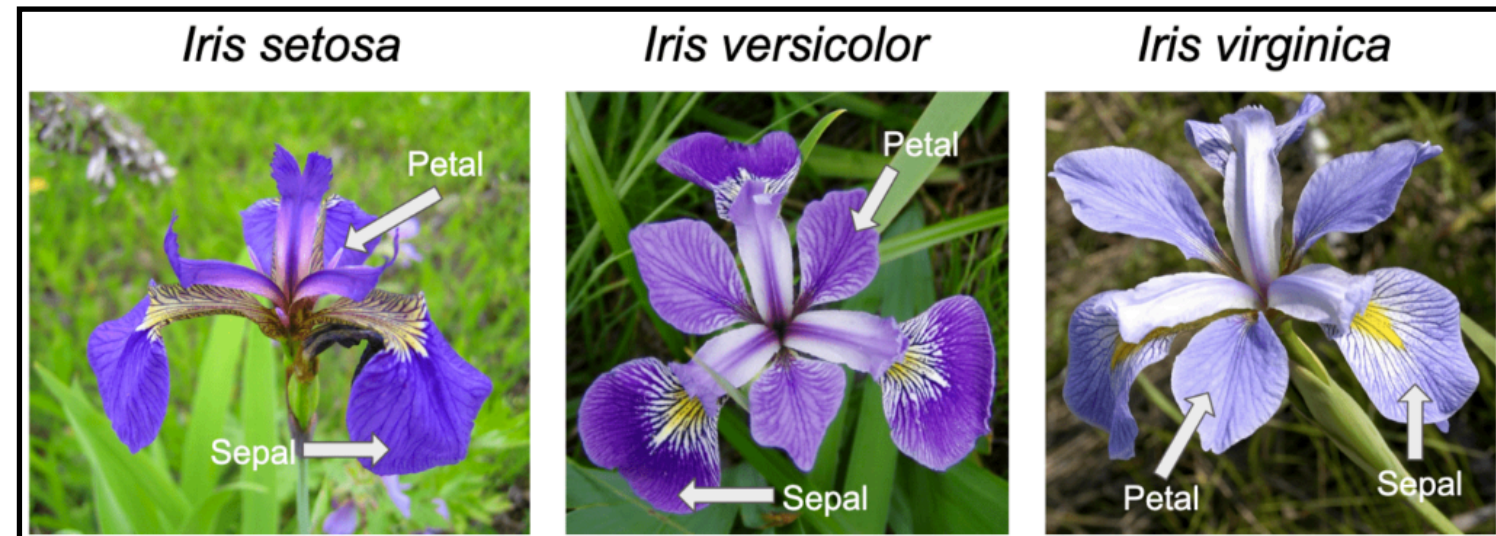
Class labels
(targets)

The diagram shows a blue Iris flower with yellow arrows indicating measurements. One arrow points vertically across the top petal, labeled 'Petal'. Another arrow points horizontally across the top petal, also labeled 'Petal'. A third arrow points diagonally across the bottom sepal, labeled 'Sepal'.

150 data points, 3 classes

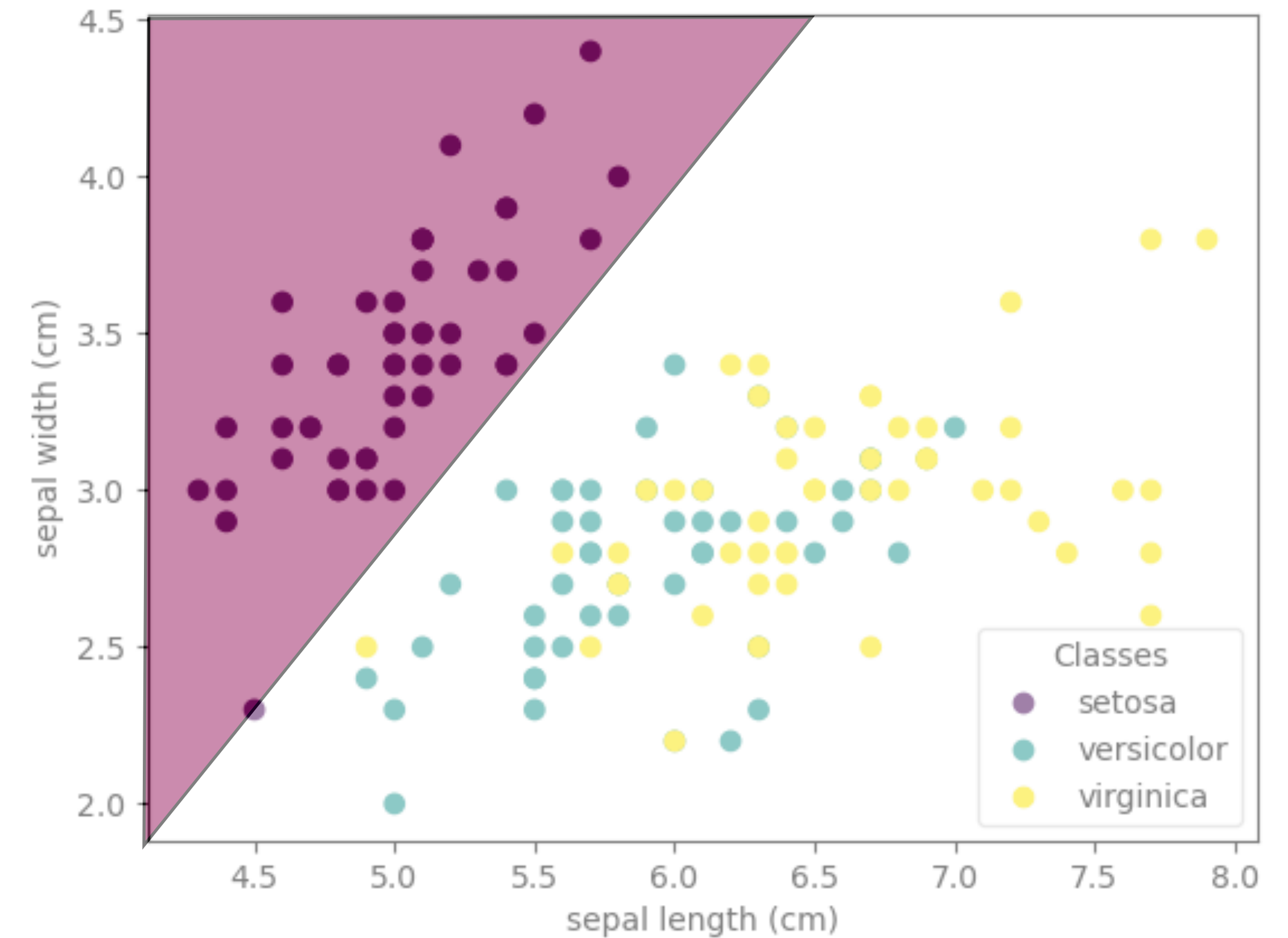
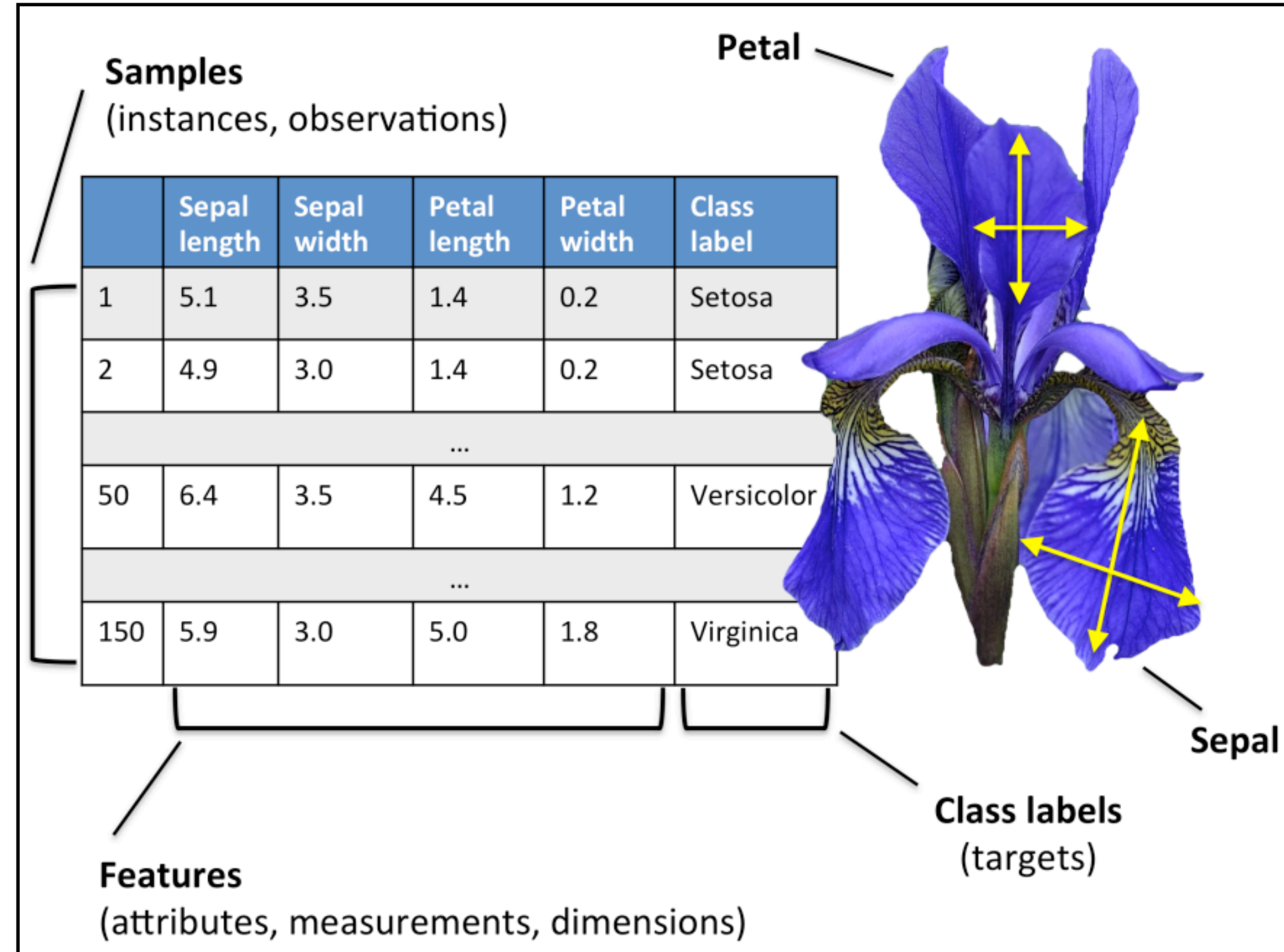
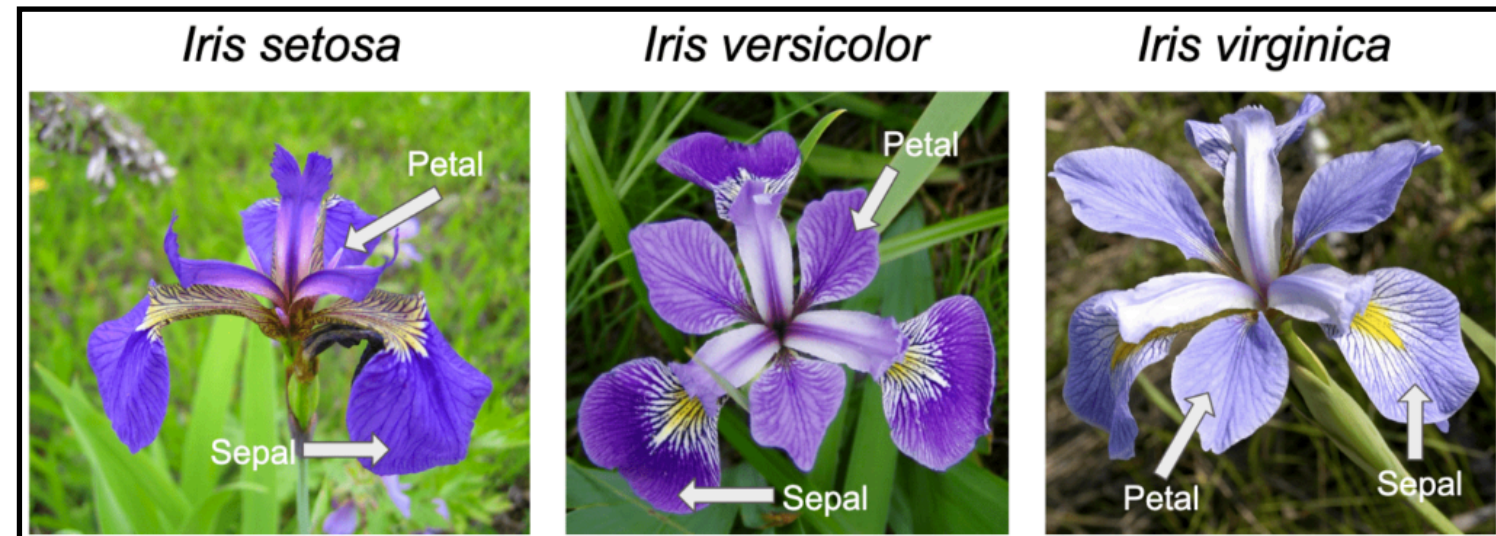
4 features

Application: Multi-Class Classification



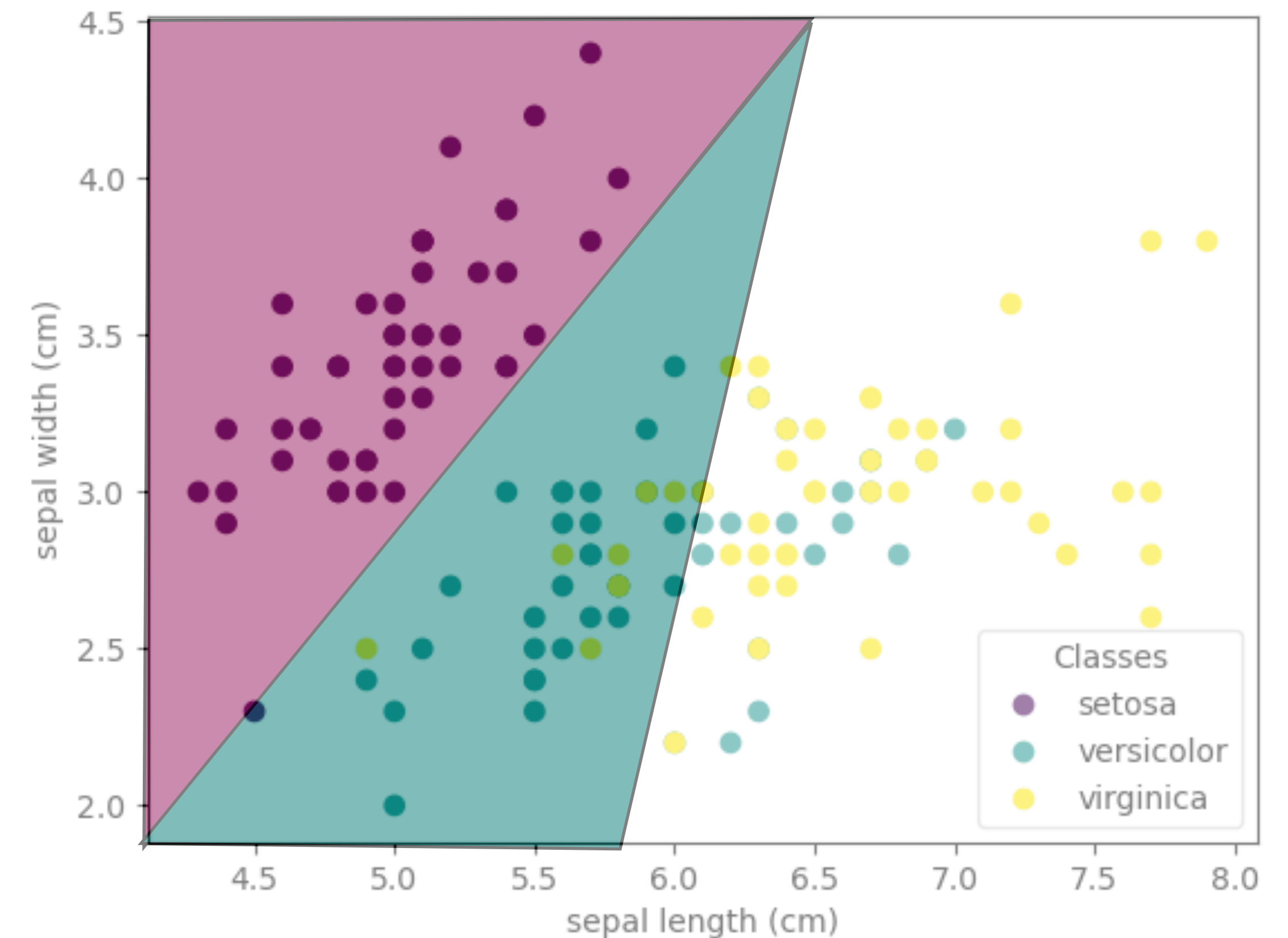
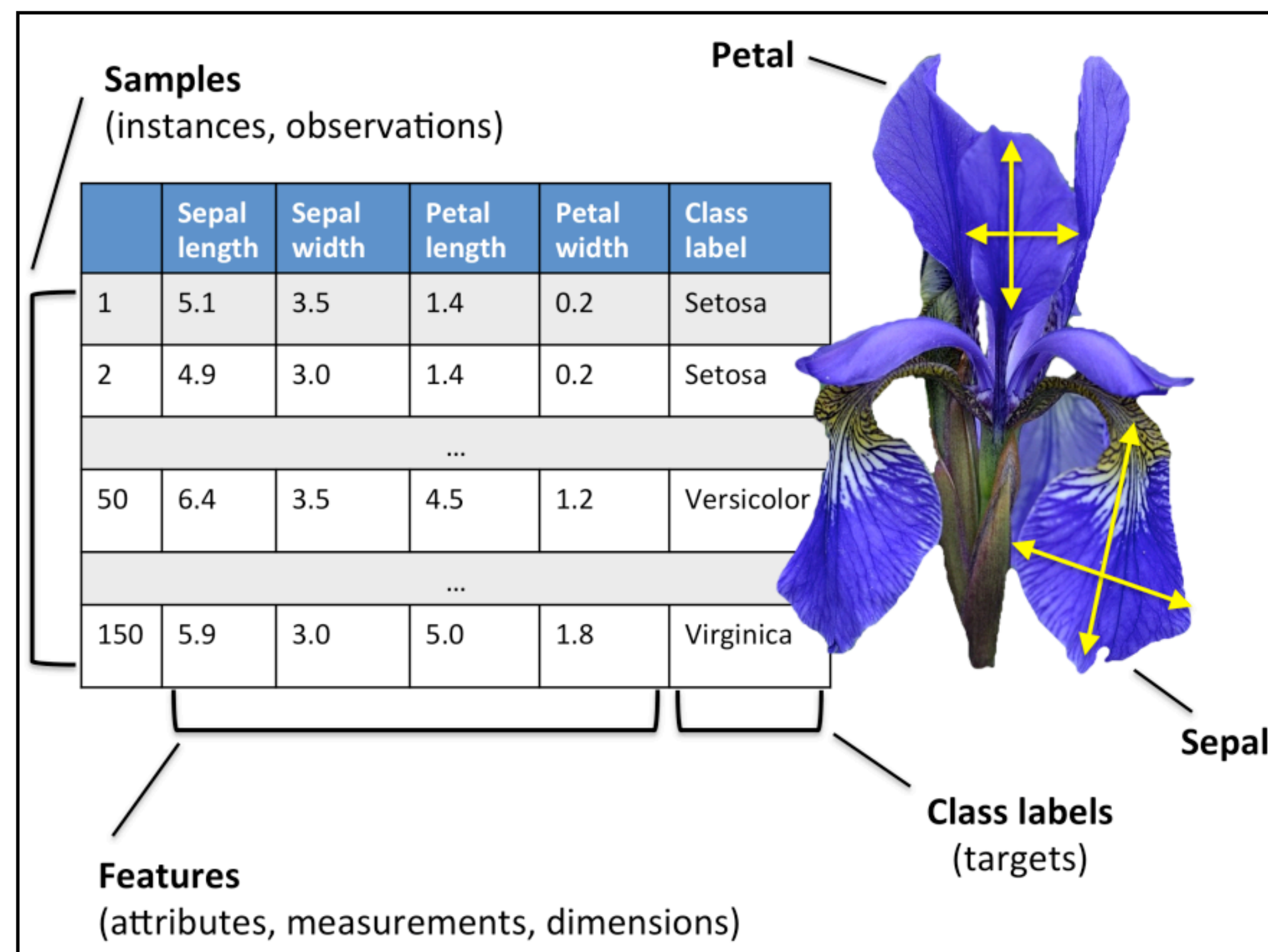
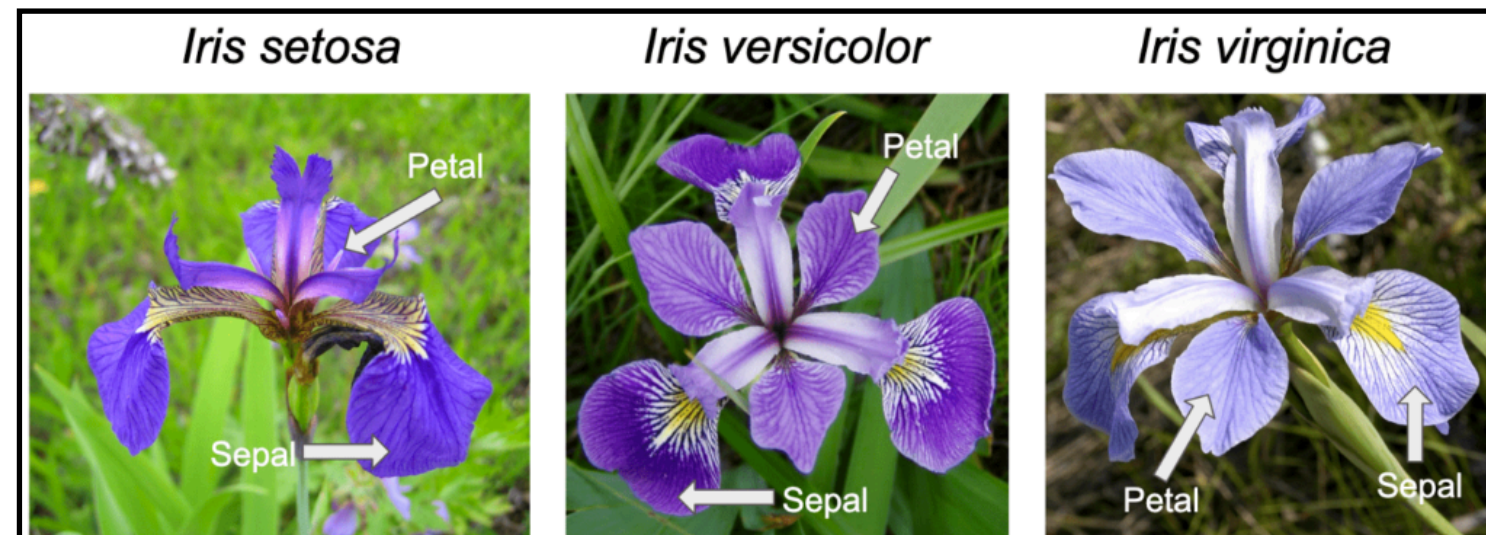
150 data points, 3 classes
4 features

Application: Multi-Class Classification



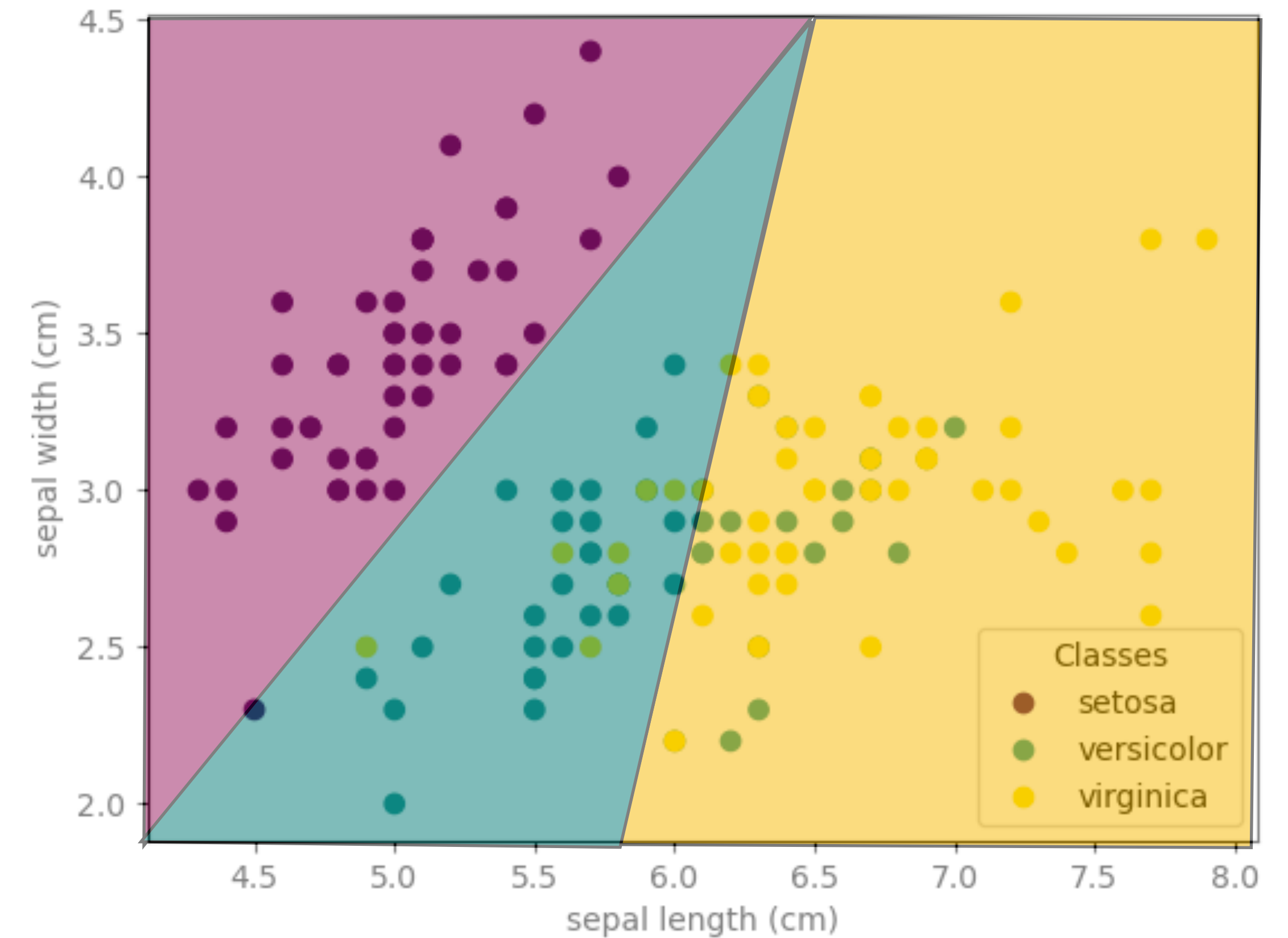
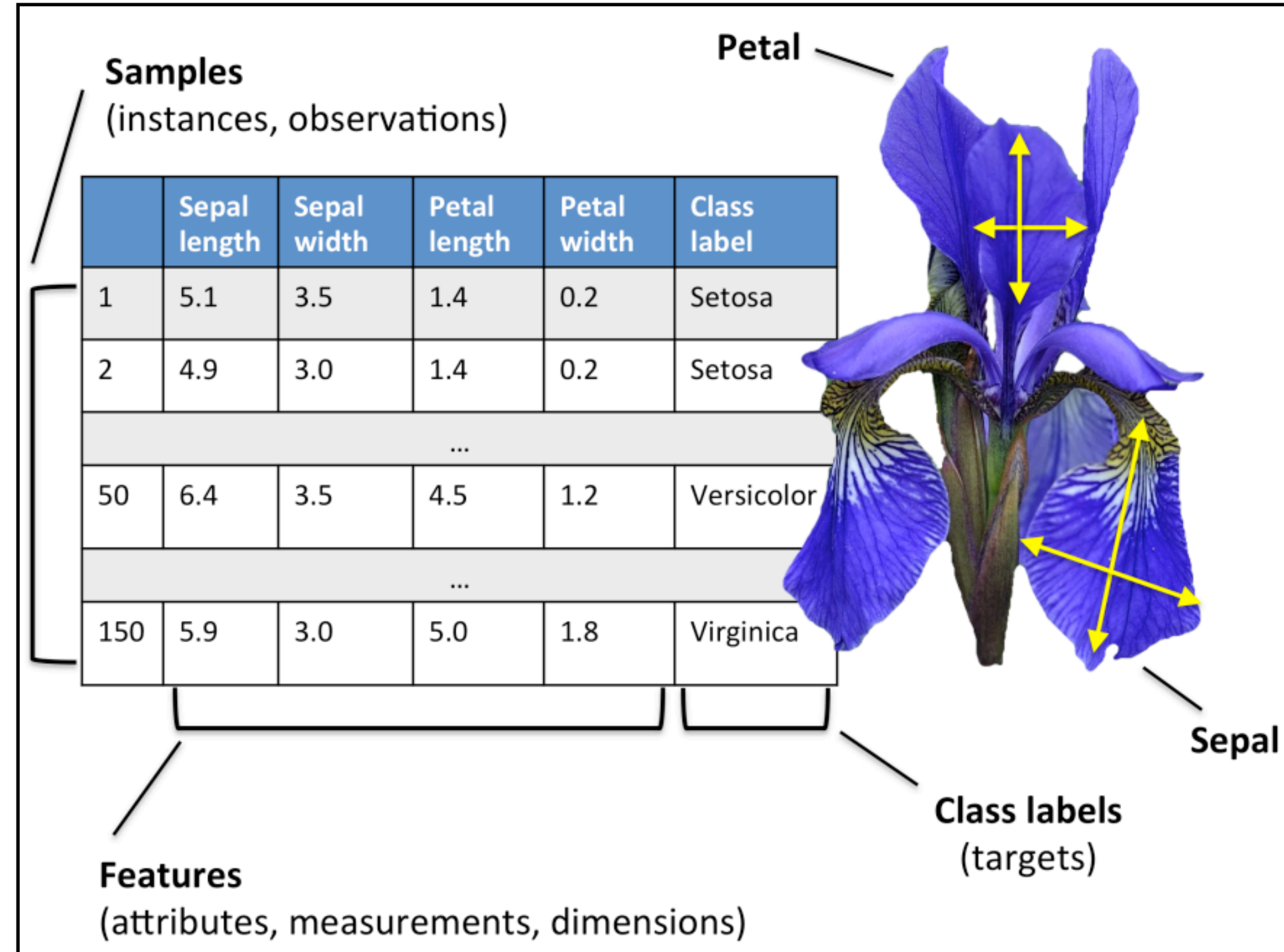
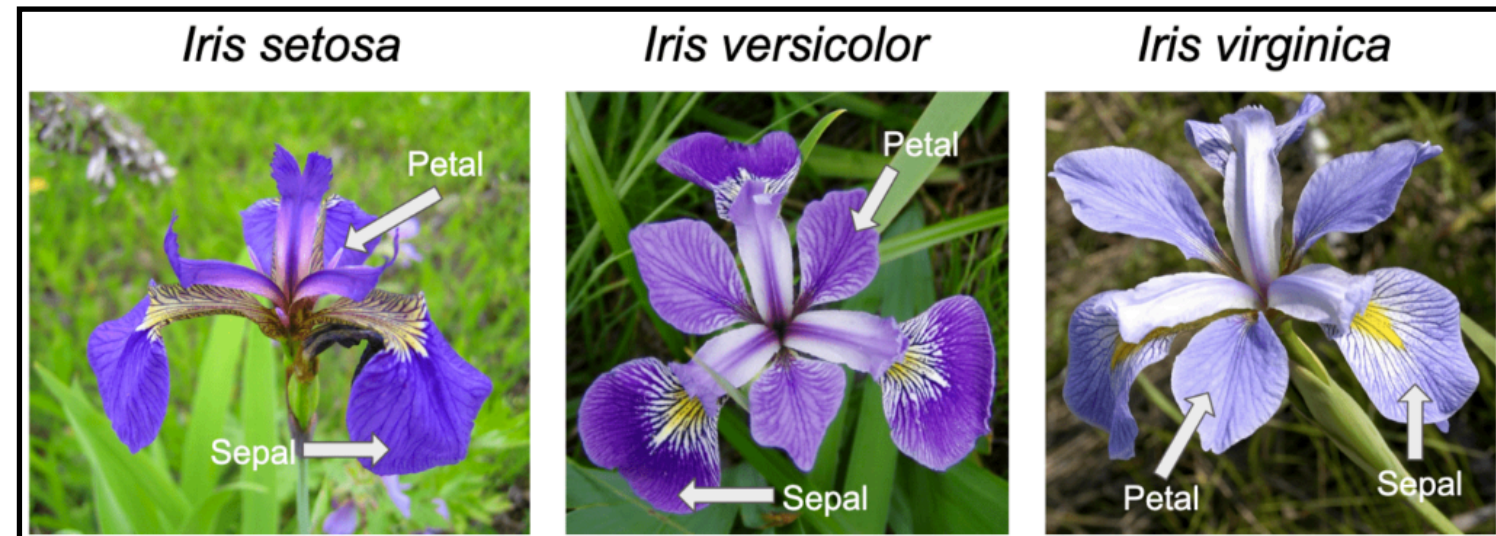
150 data points, 3 classes
4 features

Application: Multi-Class Classification



150 data points, 3 classes
4 features

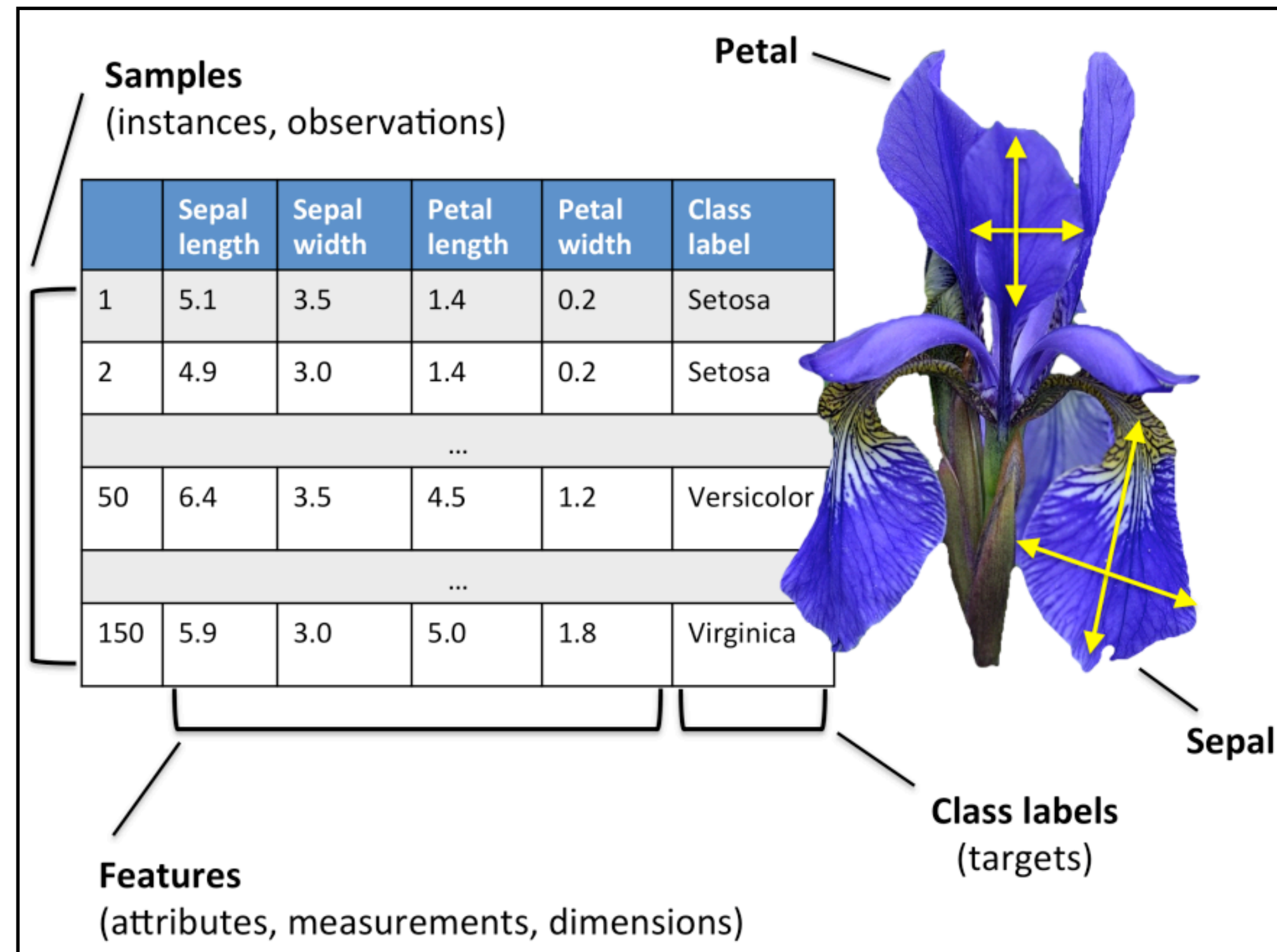
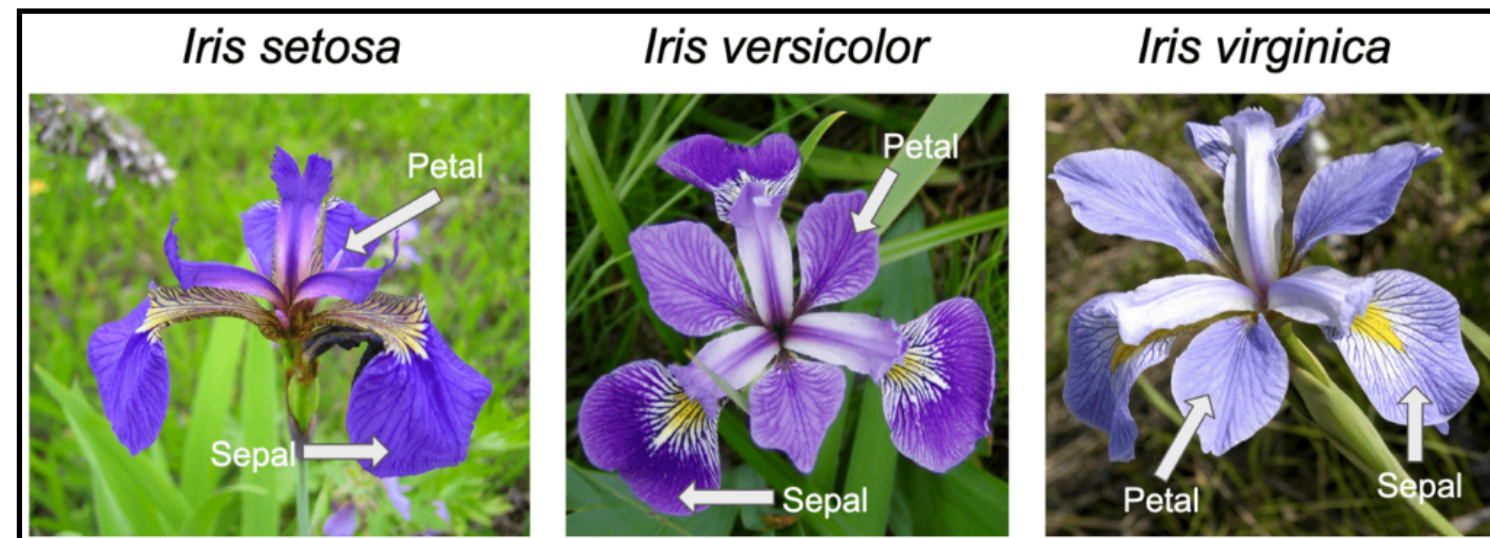
Application: Multi-Class Classification



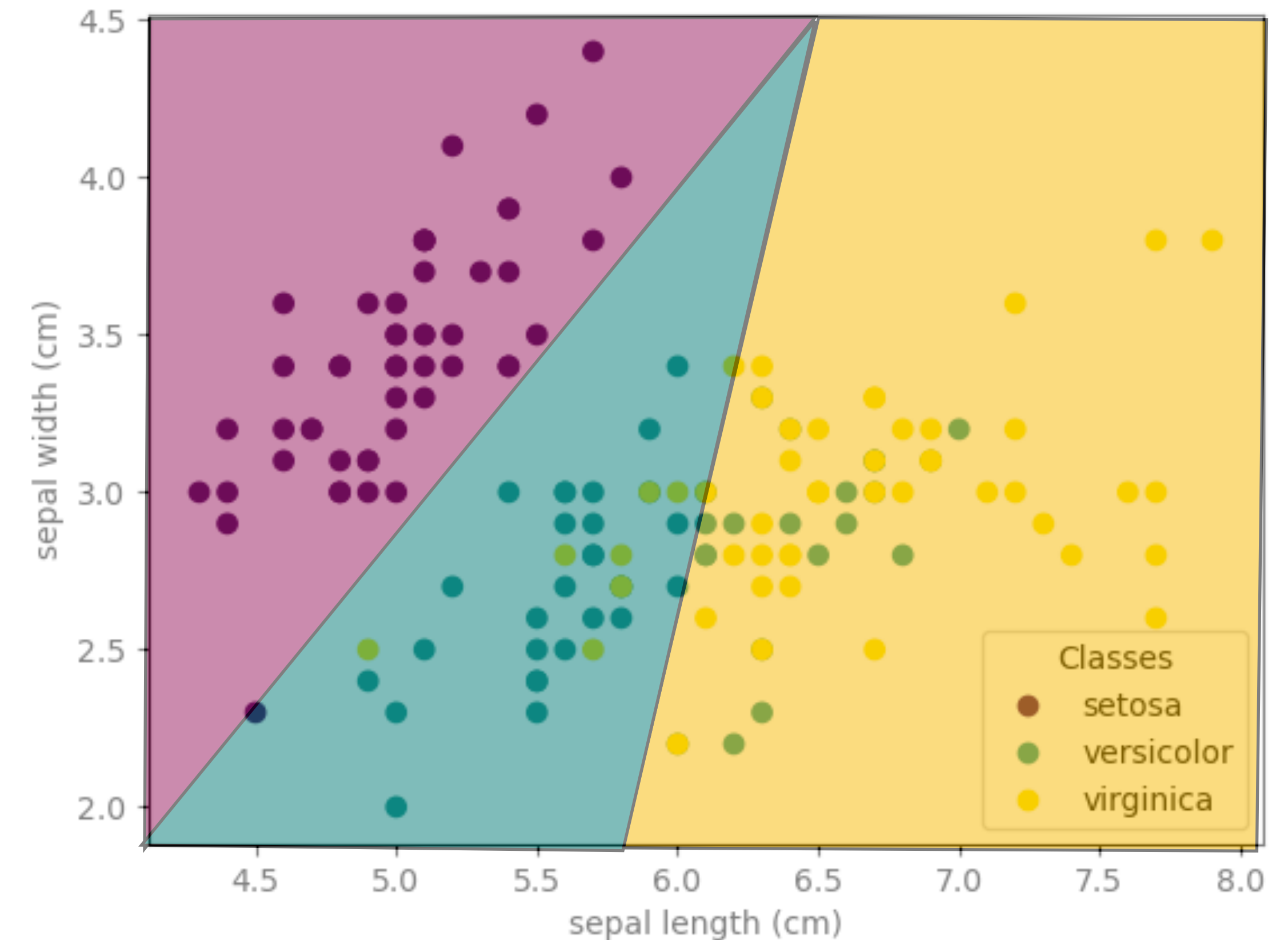
150 data points, 3 classes

4 features

Application: Multi-Class Classification

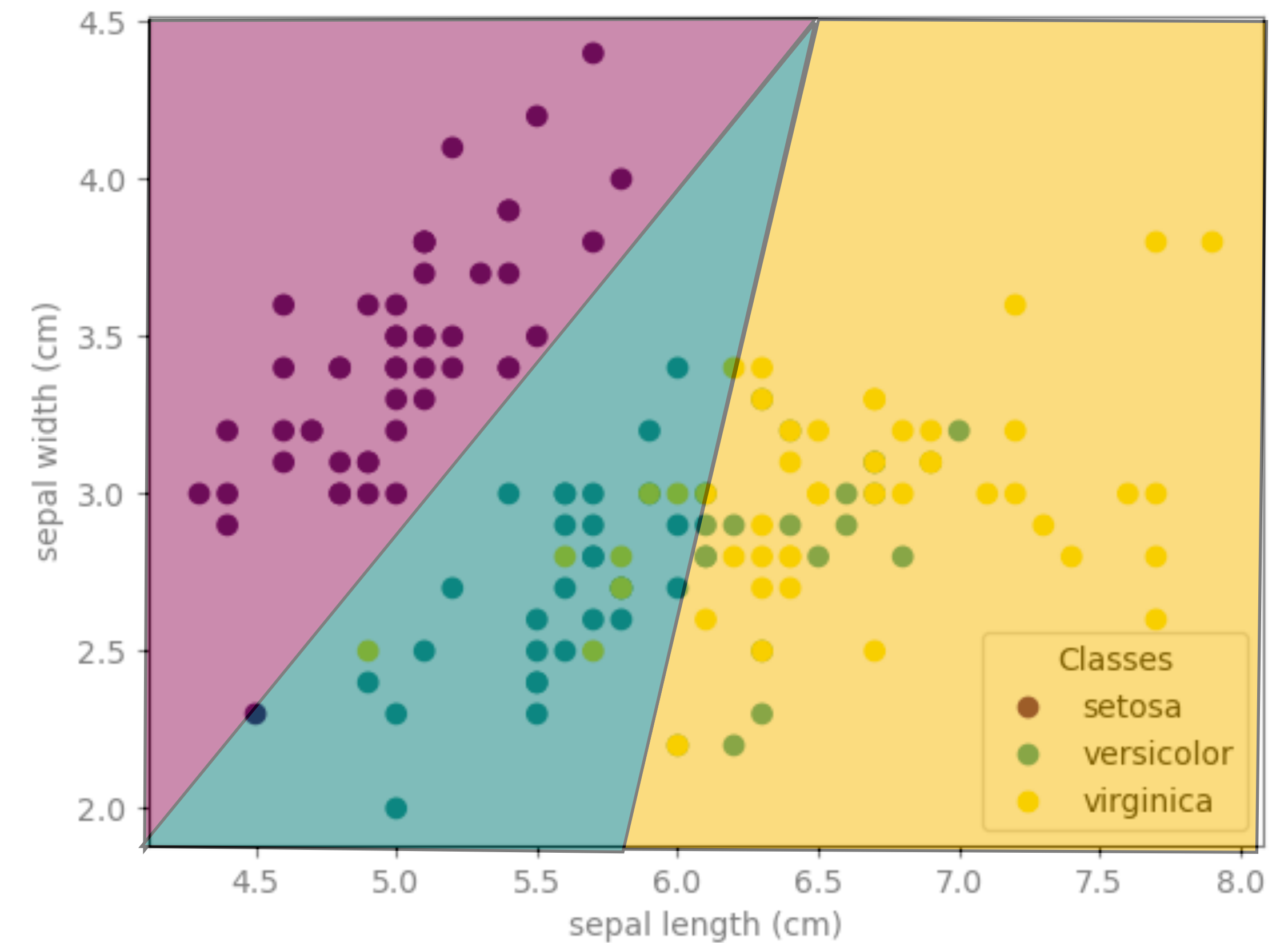
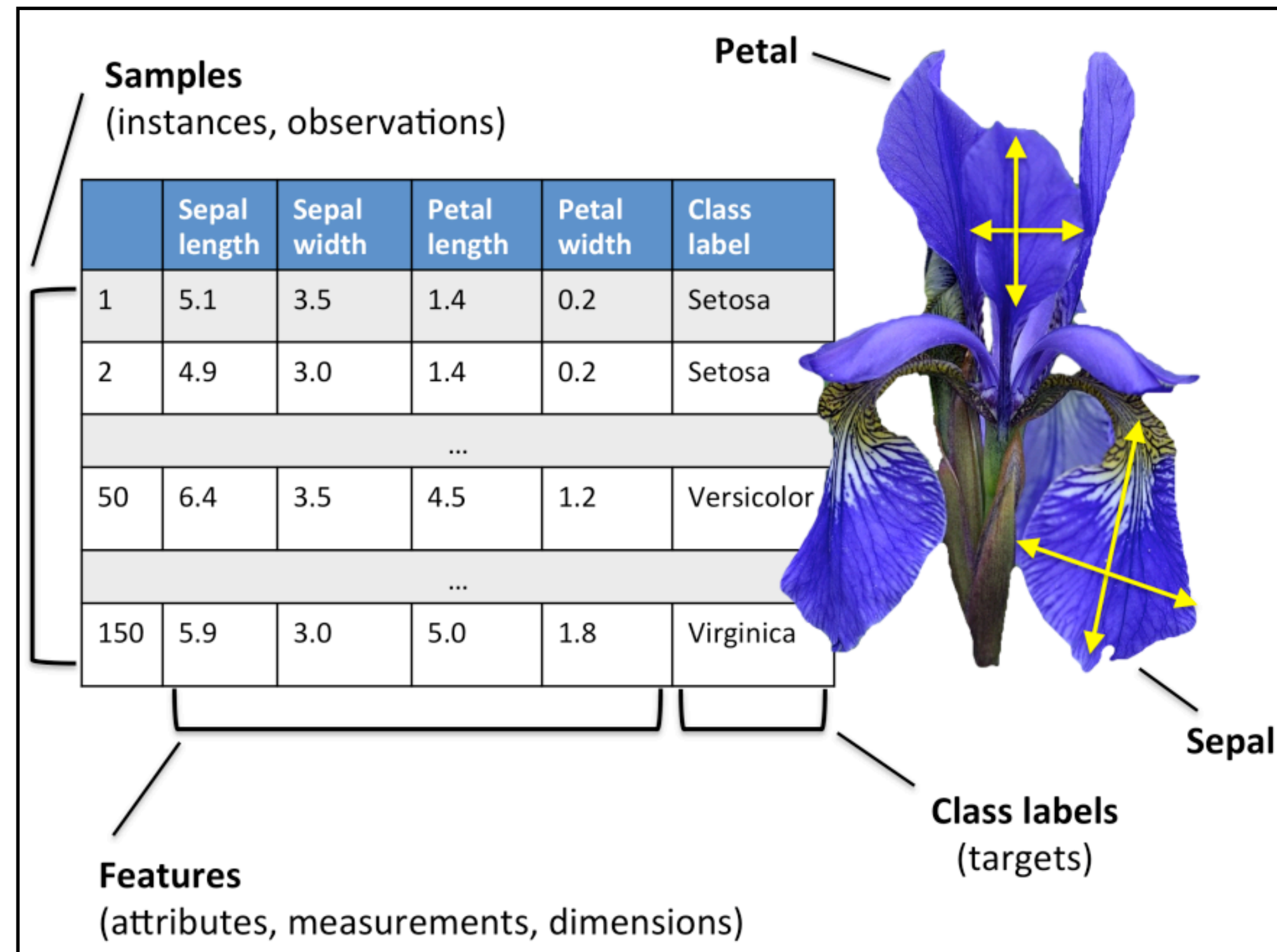
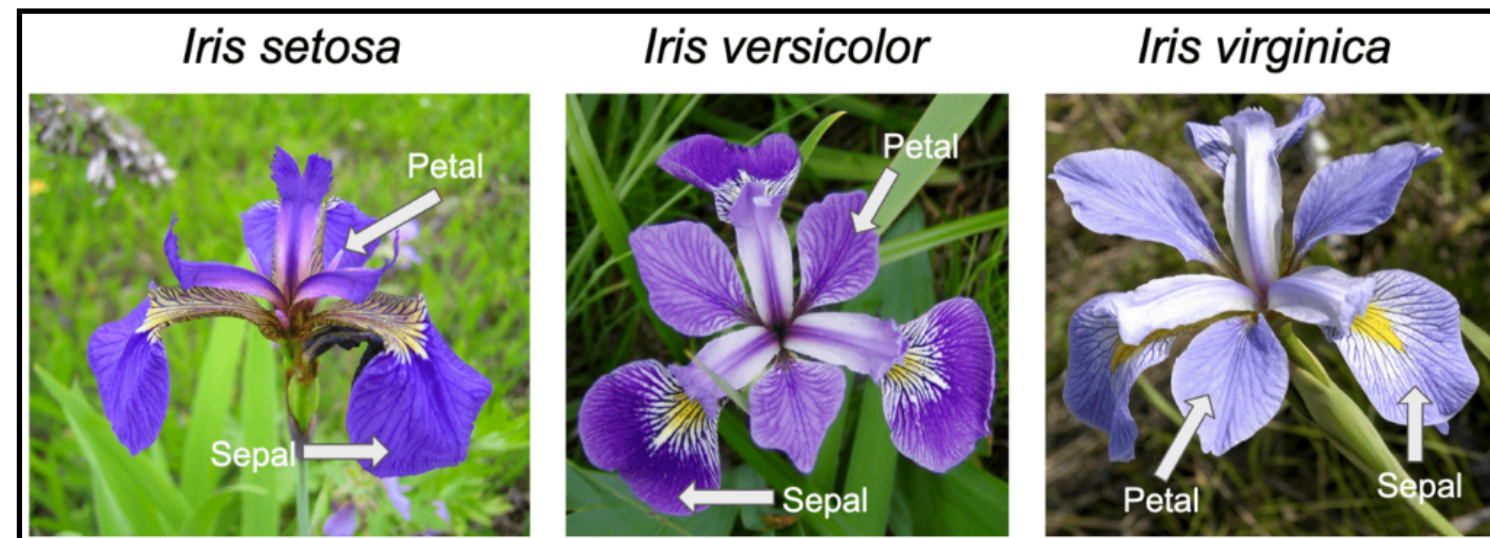


150 data points, 3 classes
4 features



Each line can be represented as

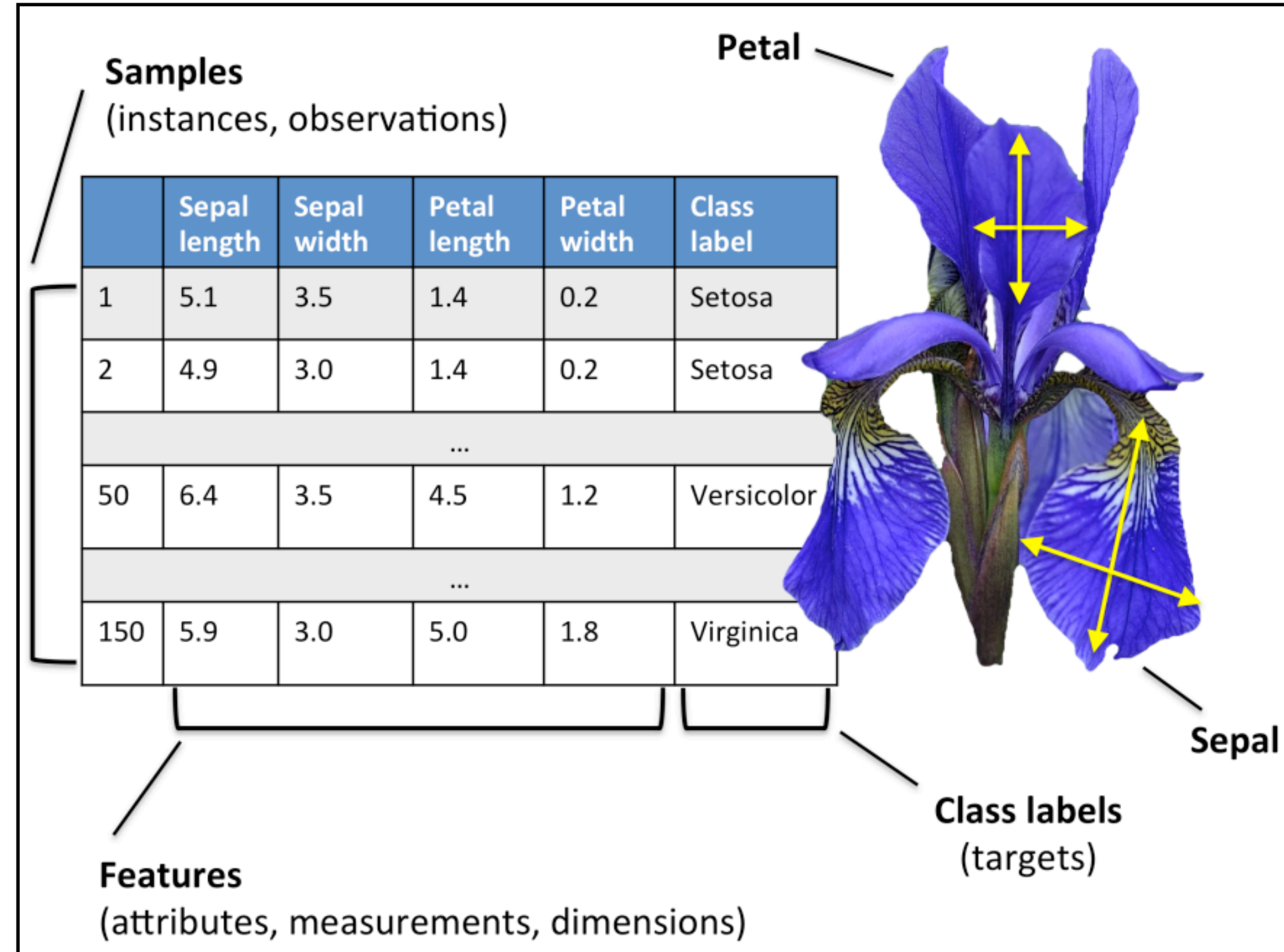
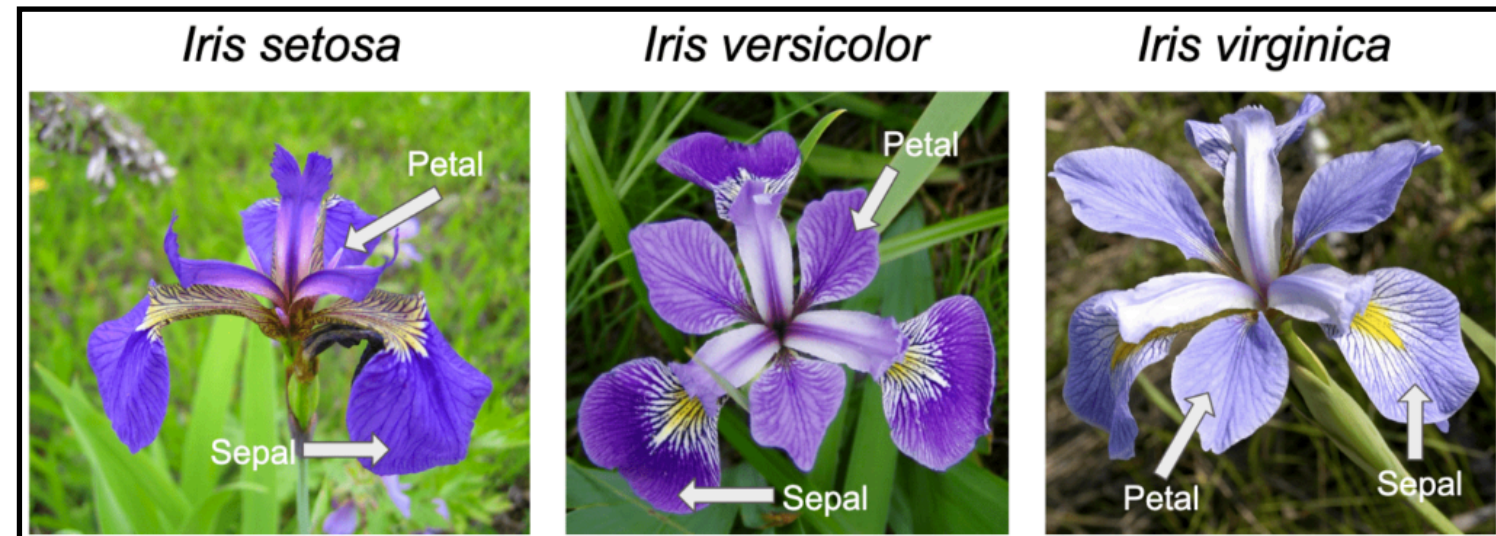
Application: Multi-Class Classification



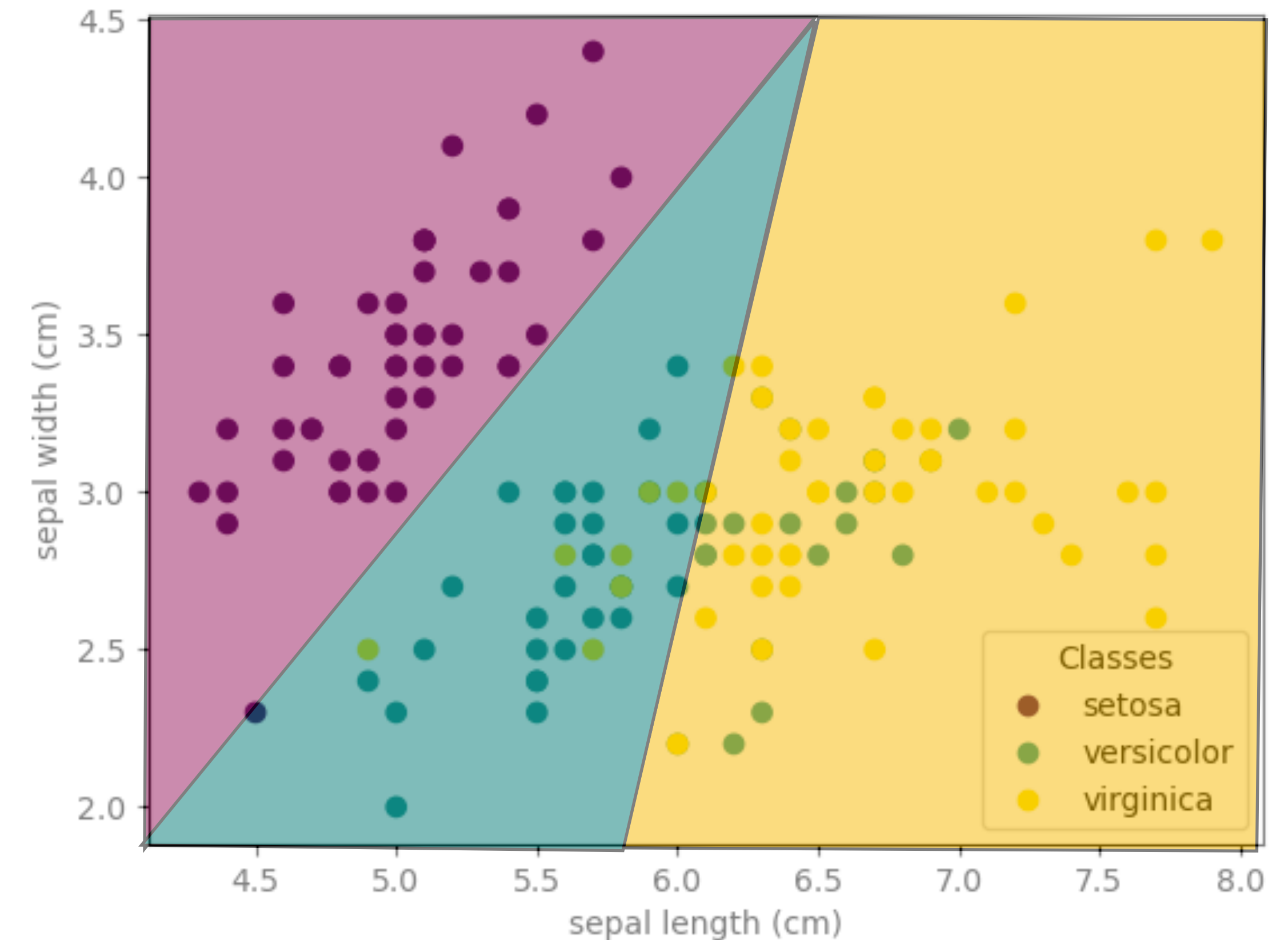
Each line can be represented as
 $w_1 \cdot (\text{sepal length}) + w_2 \cdot (\text{sepal width}) + b = 0$

150 data points, 3 classes
4 features

Application: Multi-Class Classification



150 data points, 3 classes
4 features

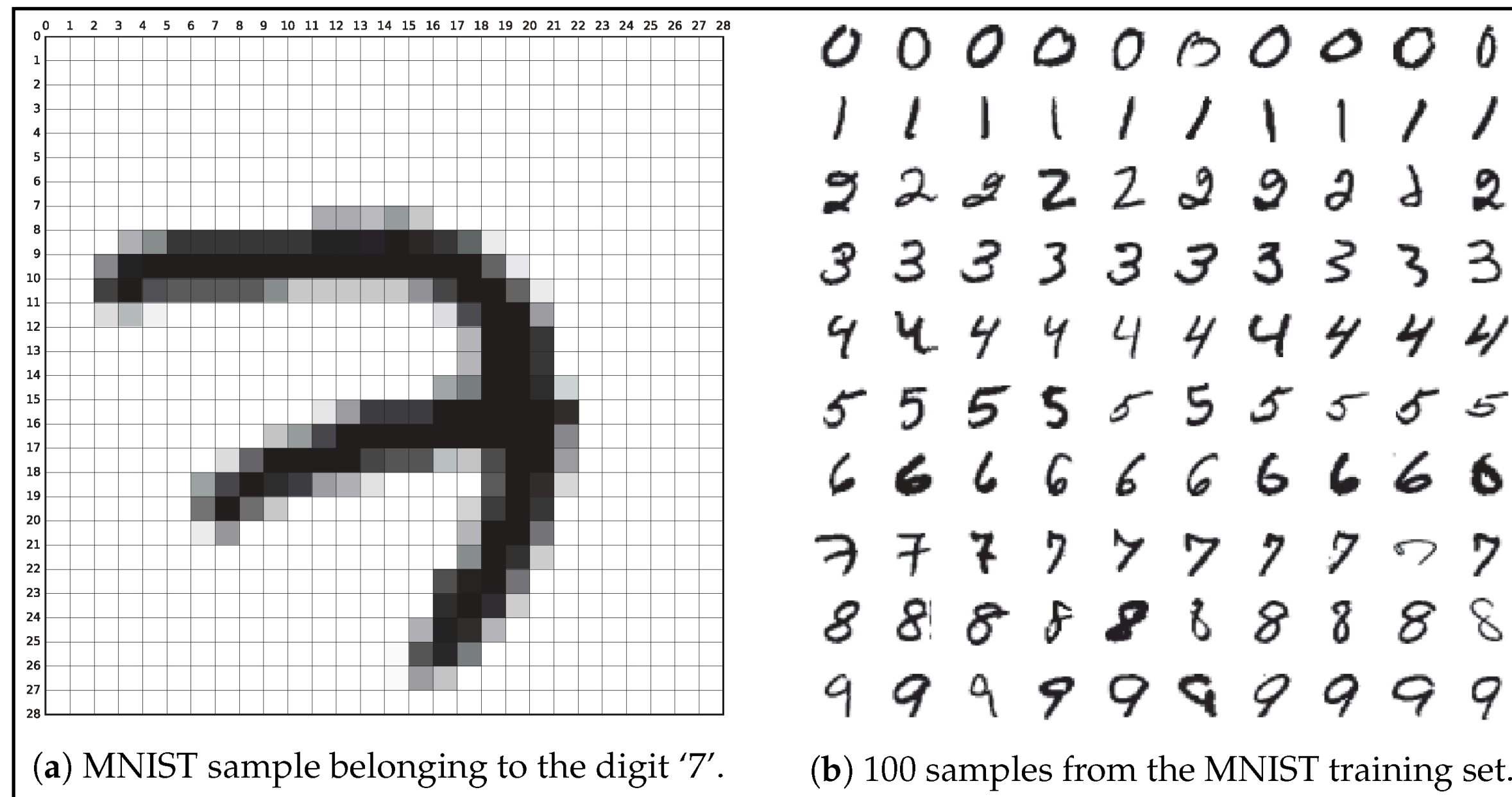


Each line can be represented as
 $w_1 \cdot (\text{sepal length}) + w_2 \cdot (\text{sepal width}) + b = 0$

Machine Learning algorithm **learns** w_1, w_2, b from the data

Application: Digit Classification (OCR)

Application: Digit Classification (OCR)

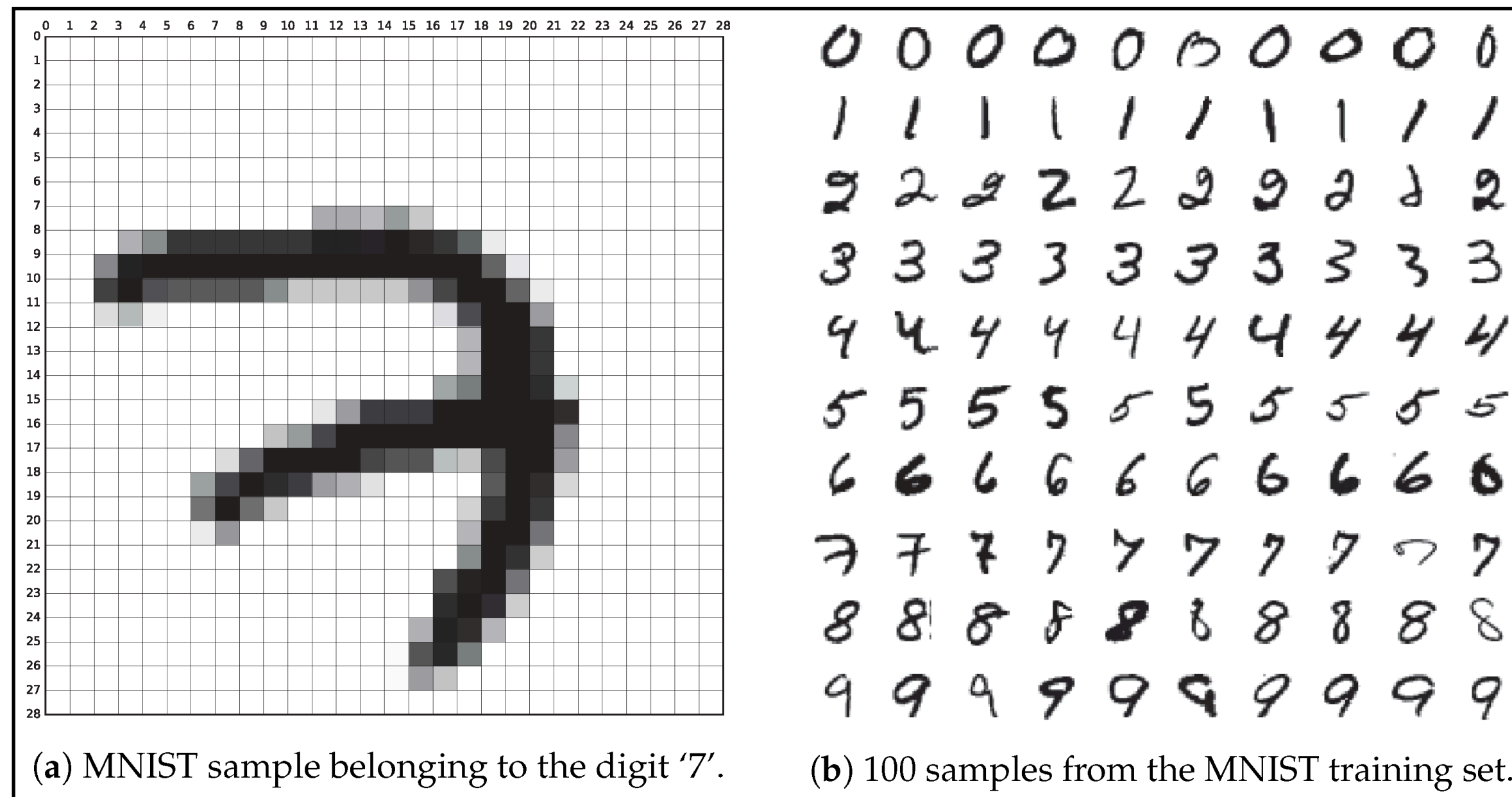


MNIST is a digit dataset:

50,000 images,

28 x 28 features (728) per image

Application: Digit Classification (OCR)



Linear classifiers i.e.

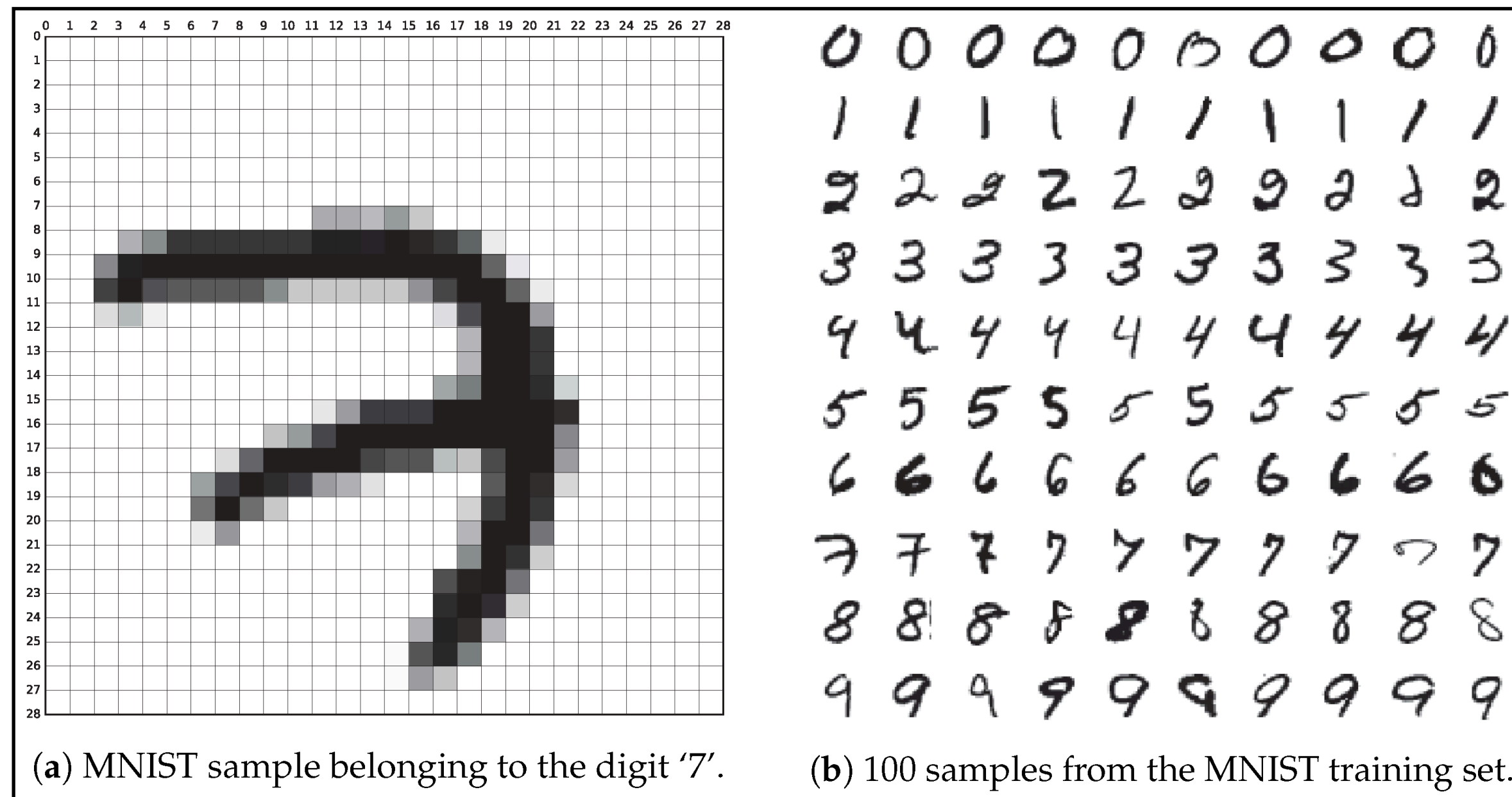
$w_1x_1 + w_2x_2 + \dots + w_{728}x_{728} + b = 0$
are not sufficient here.

MNIST is a digit dataset:

50,000 images,

28 x 28 features (728) per image

Application: Digit Classification (OCR)



Linear classifiers i.e.

$w_1x_1 + w_2x_2 + \dots + w_{728}x_{728} + b = 0$
are not sufficient here.

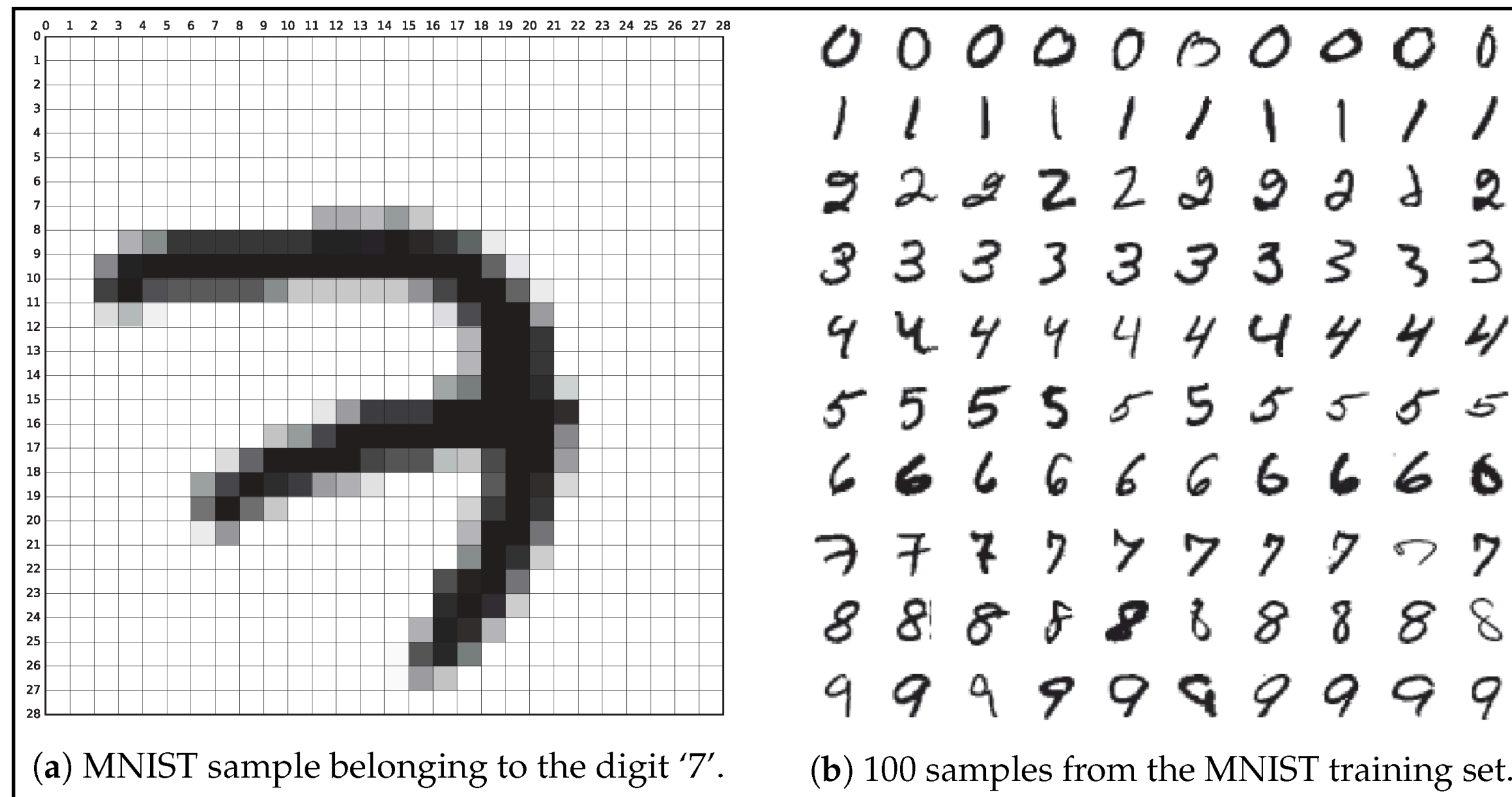
We use Artificial Neural Networks

MNIST is a digit dataset:

50,000 images,

28 x 28 features (728) per image

Application: Digit Classification (OCR)

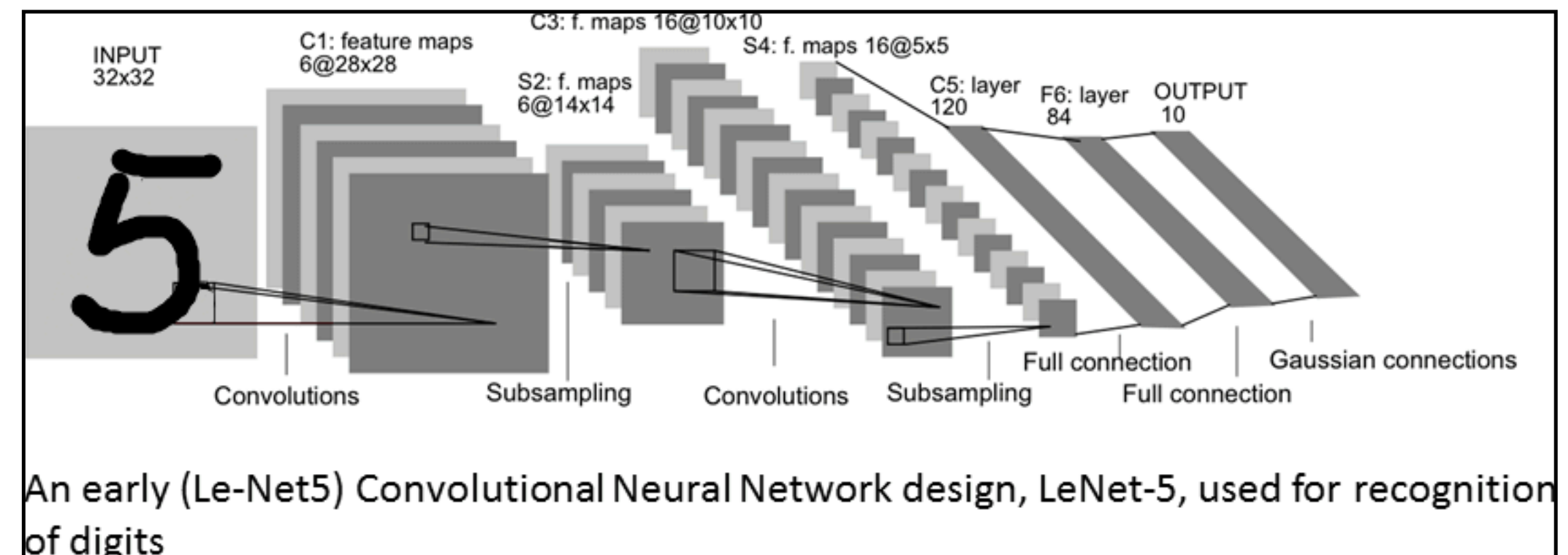


MNIST is a digit dataset:
50,000 images,
28 x 28 features (728) per image

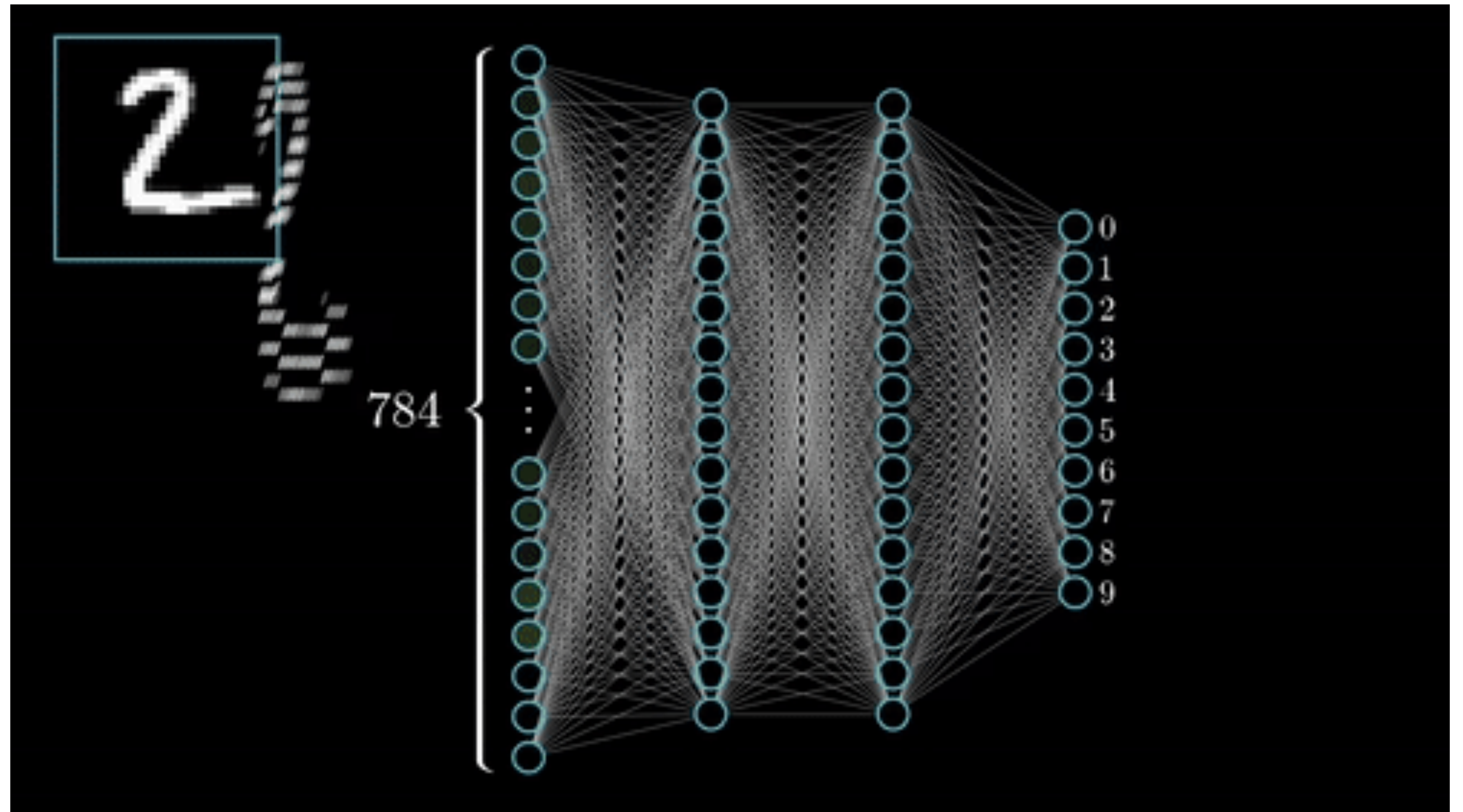
Linear classifiers i.e.

$w_1x_1 + w_2x_2 + \dots + w_{728}x_{728} + b = 0$
are not sufficient here.

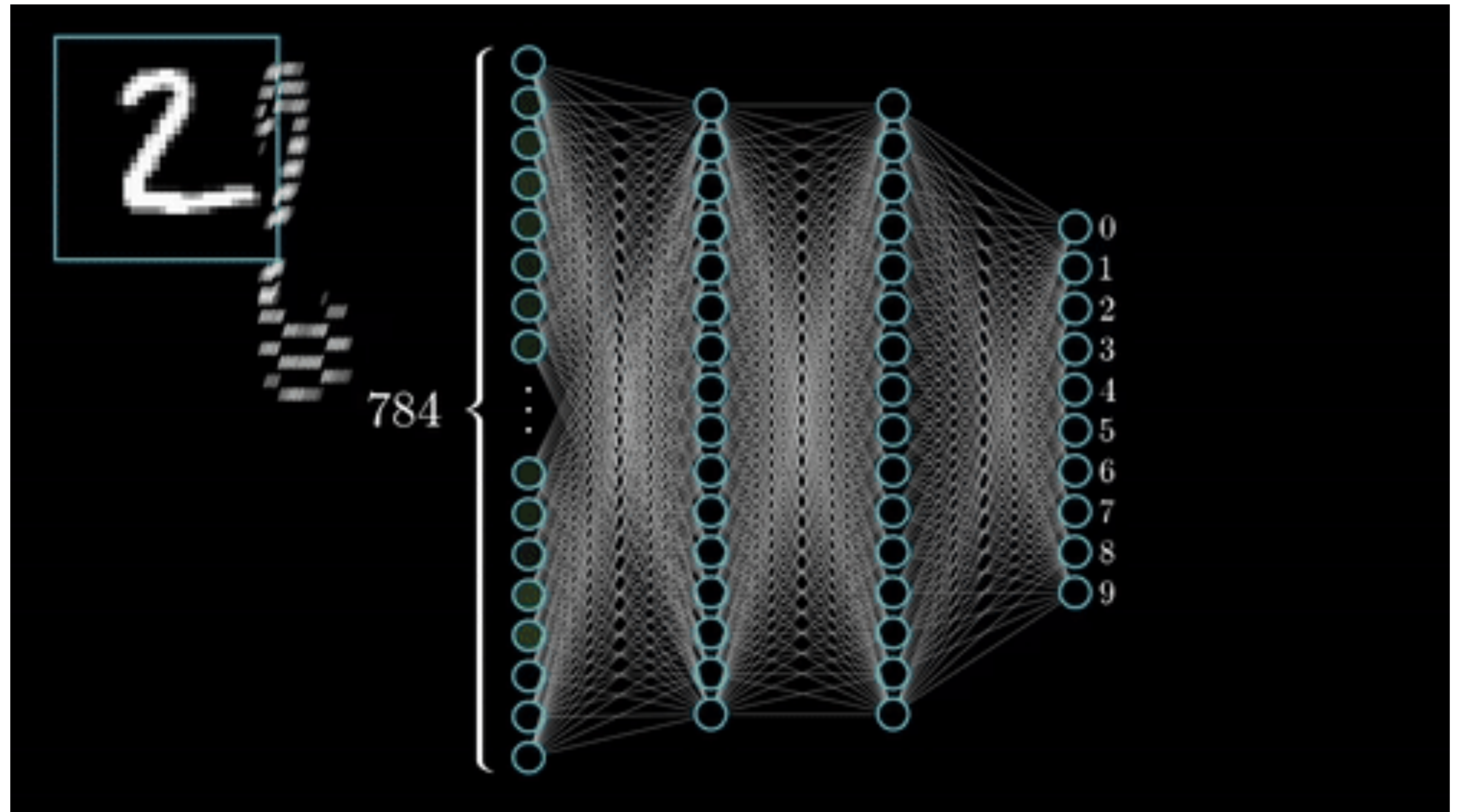
We use Artificial Neural Networks



Neural Networks: Backbone of modern ML

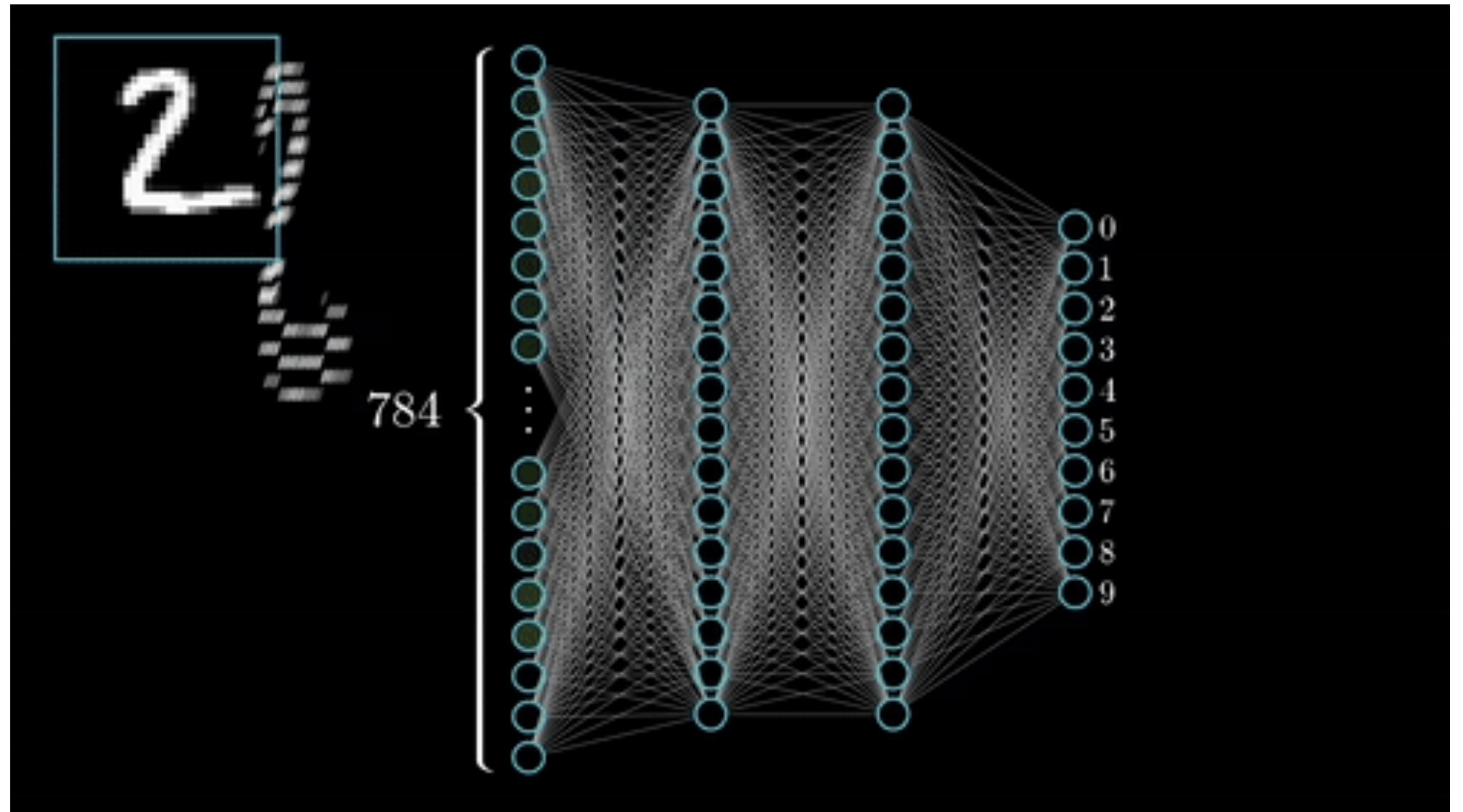


Neural Networks: Backbone of modern ML



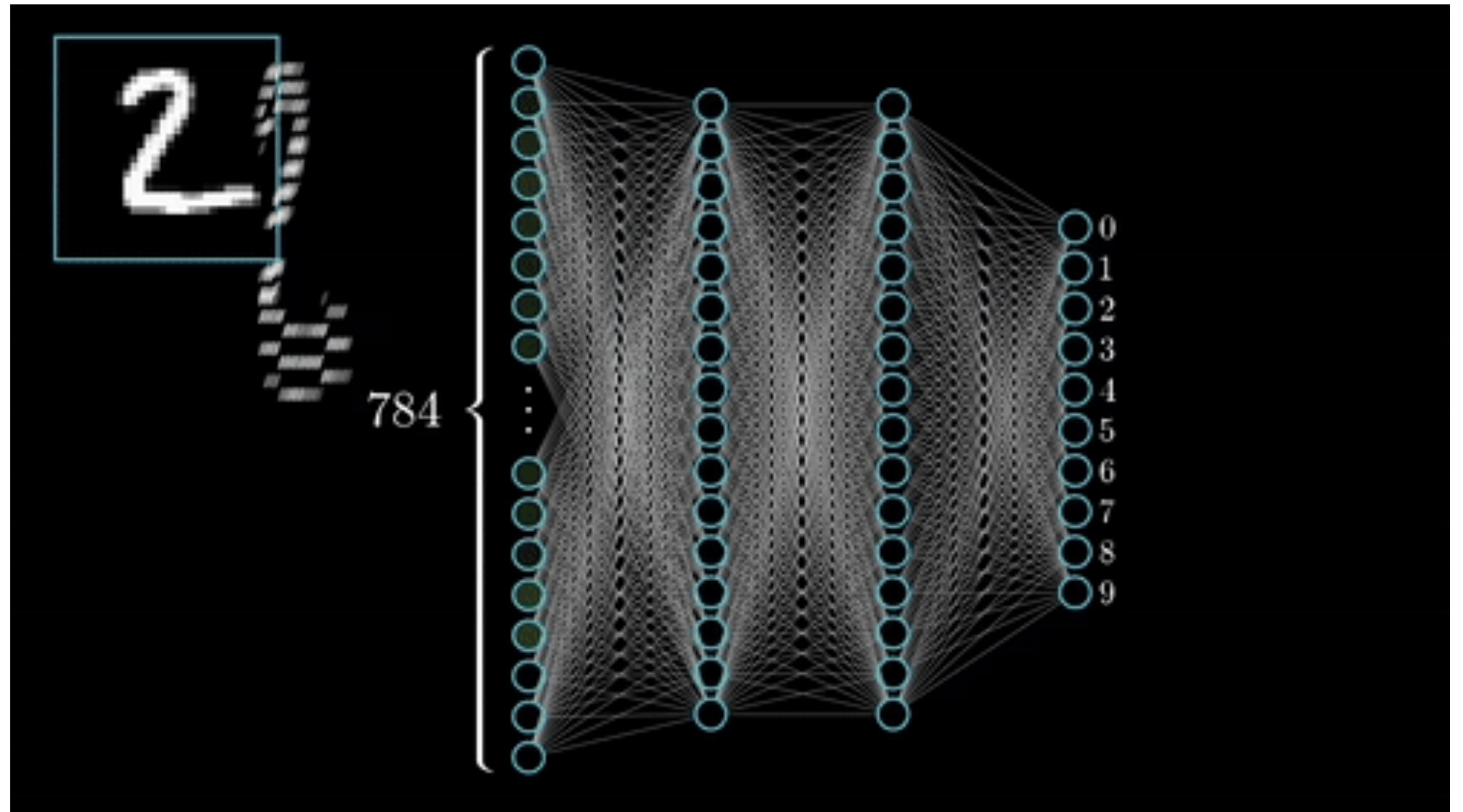
Neural Networks: Backbone of modern ML

Artificial Neural Networks are



Neural Networks: Backbone of modern ML

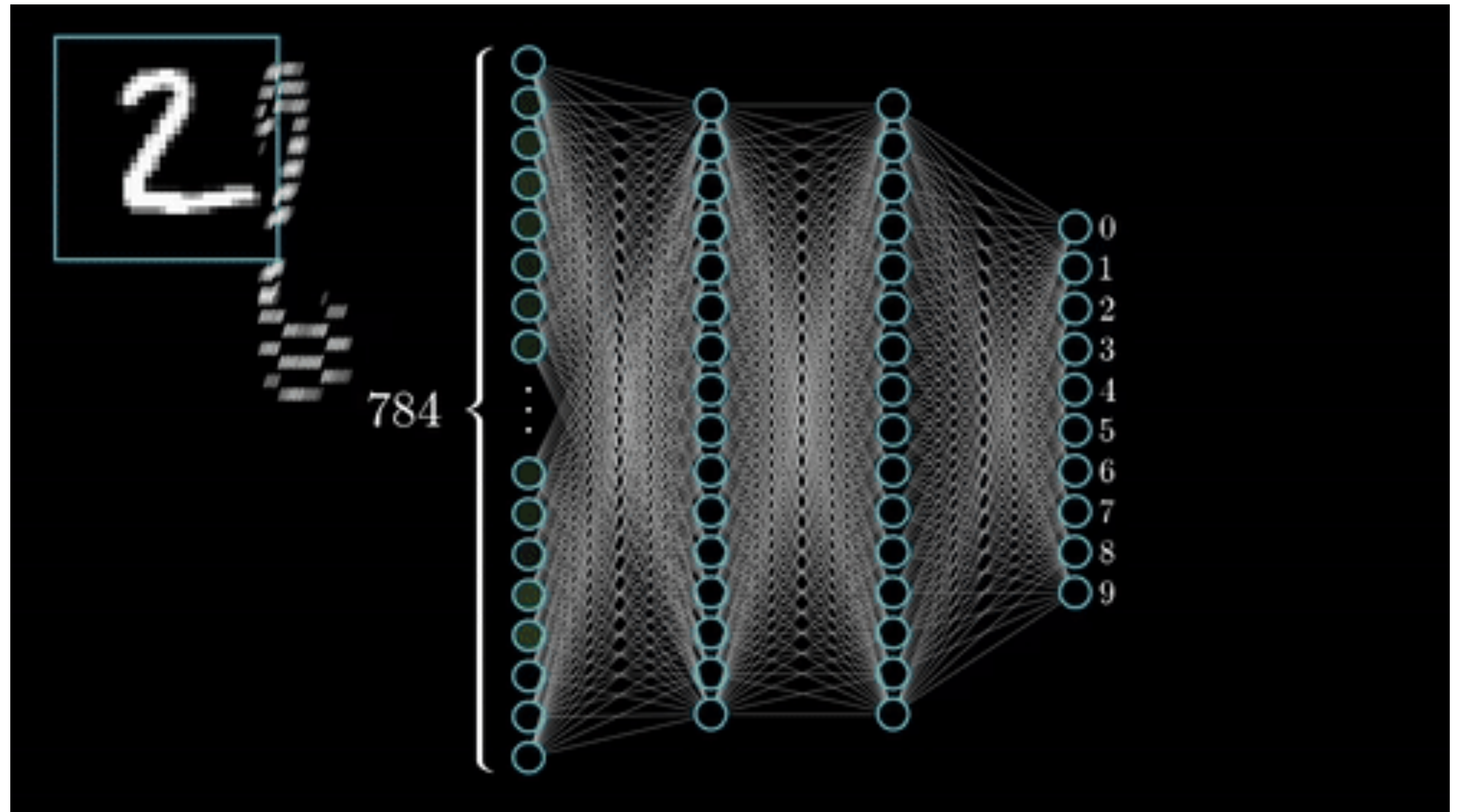
Artificial Neural Networks are
highly expressive



Neural Networks: Backbone of modern ML

Artificial Neural Networks are
highly expressive

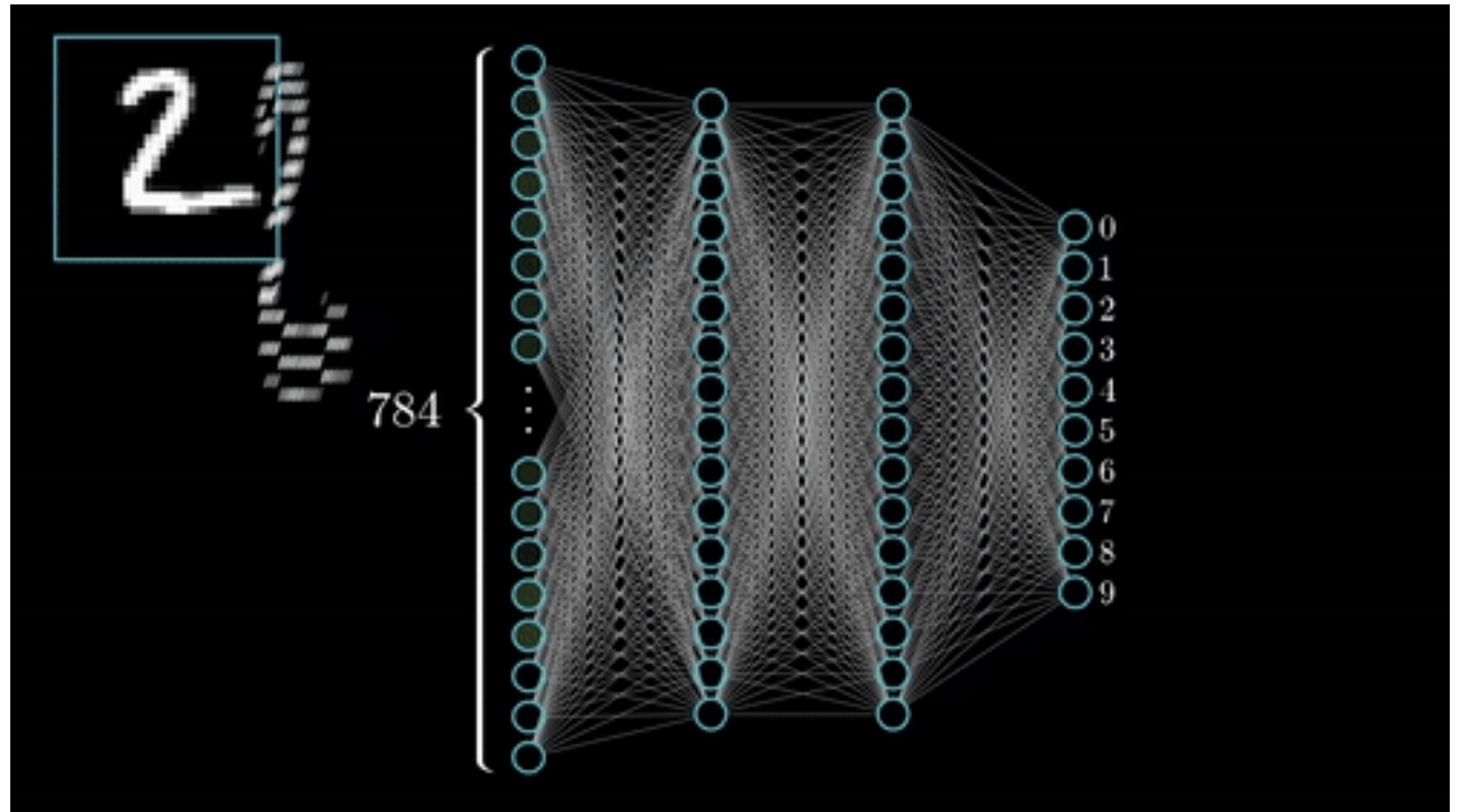
But they have



Neural Networks: Backbone of modern ML

Artificial Neural Networks are
highly expressive

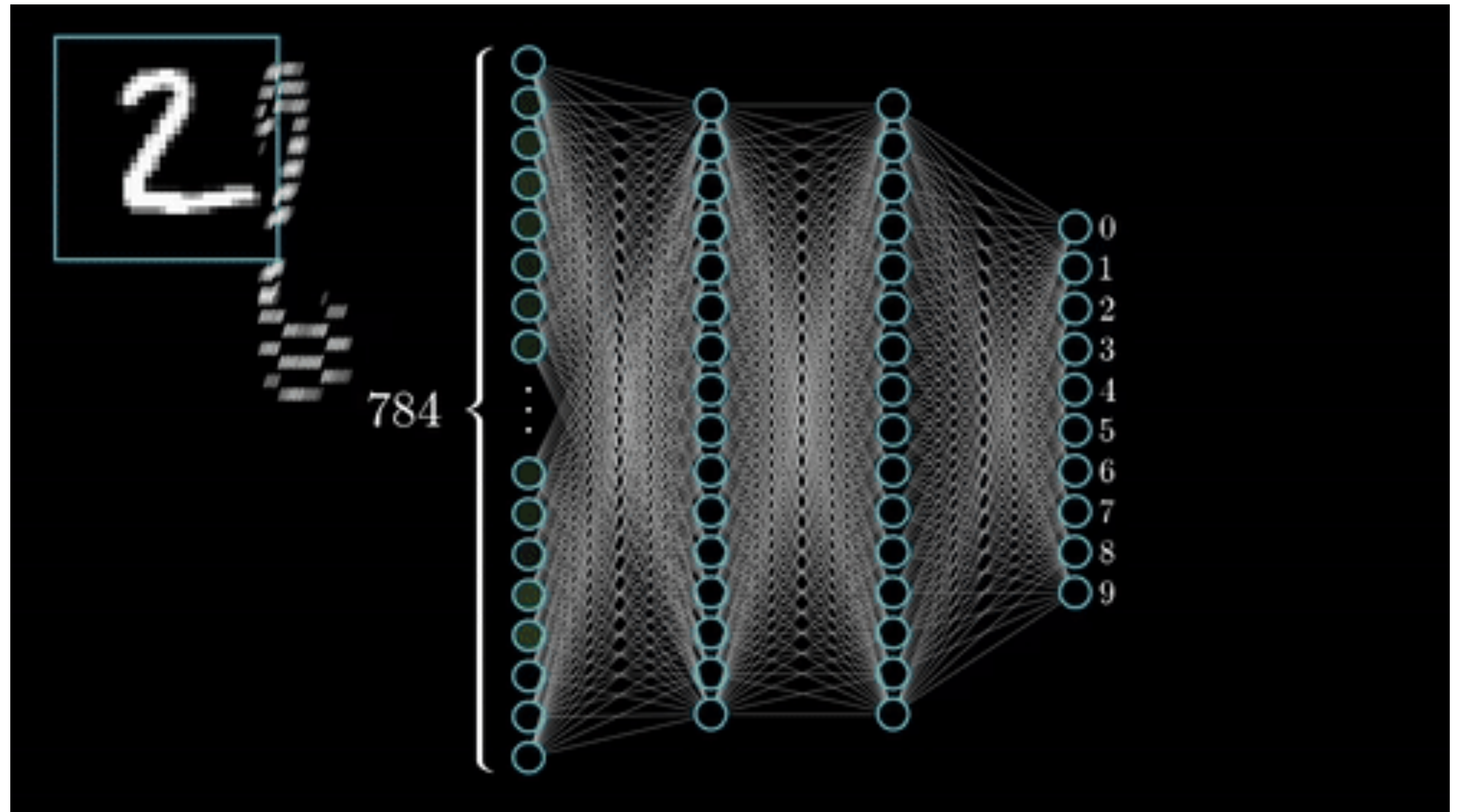
But they have
Large number of parameters



Neural Networks: Backbone of modern ML

Artificial Neural Networks are **highly expressive**

But they have
Large number of parameters

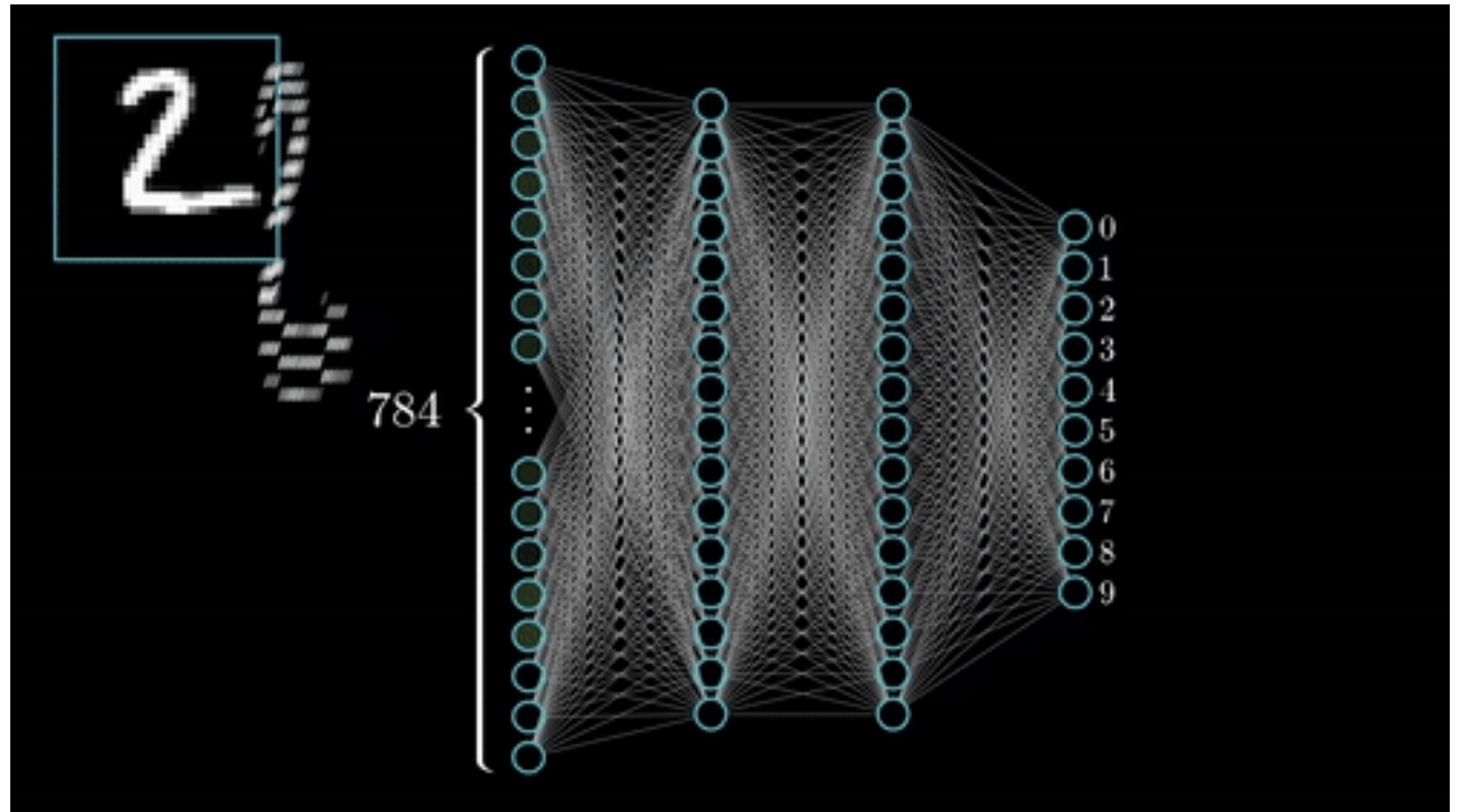


Rumors claim that GPT-4 has 1.76 trillion parameters, which was first estimated by the speed it was running and by [George Hotz](#).^[12]

Neural Networks: Backbone of modern ML

Artificial Neural Networks are **highly expressive**

But they have
Large number of parameters



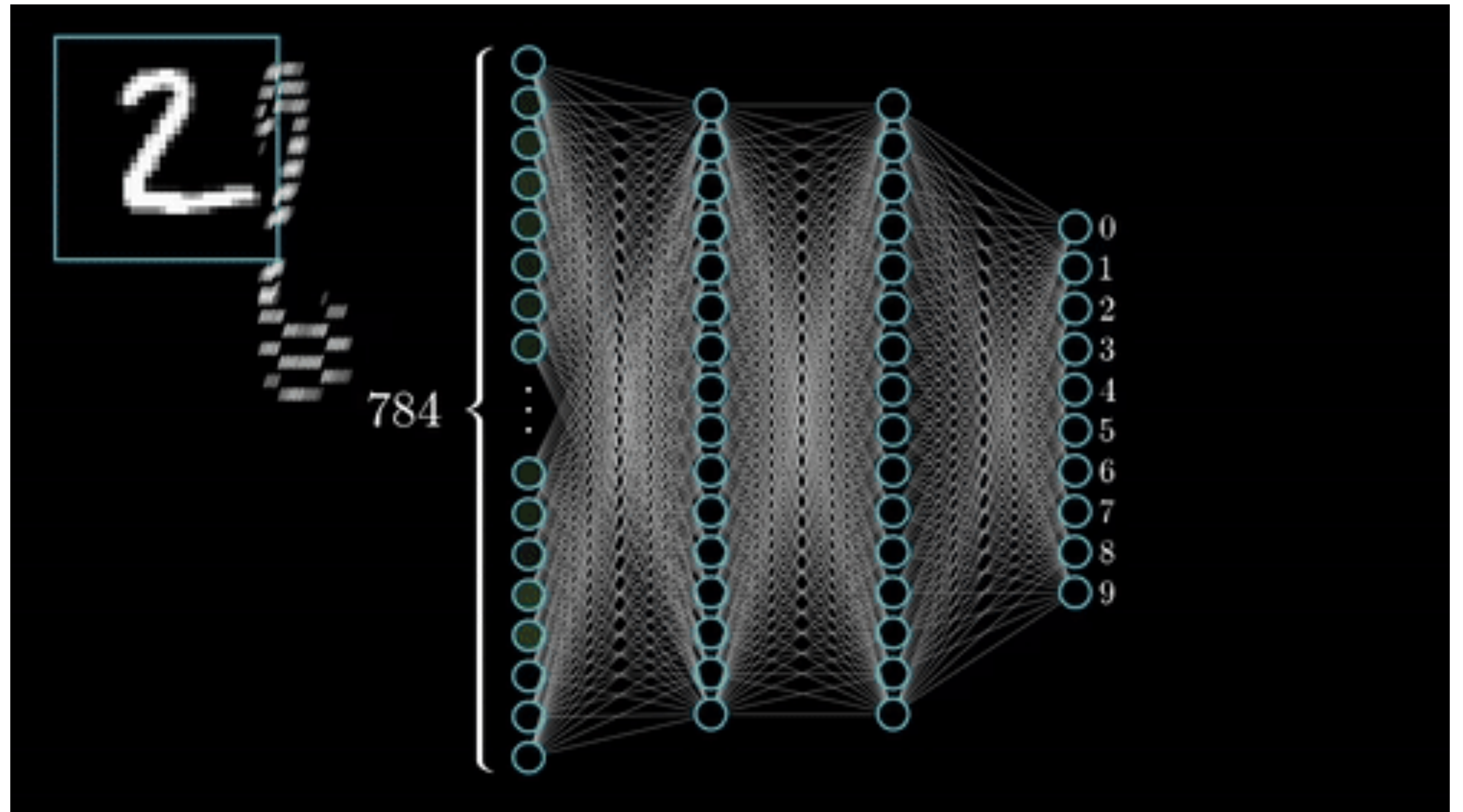
Rumors claim that [GPT-4 has 1.76 trillion parameters](#), which was first estimated by the speed it was running and by [George Hotz](#).^[12]

Neural Networks: Backbone of modern ML

Artificial Neural Networks are **highly expressive**

But they have
Large number of parameters

Each edge is a parameter.



Rumors claim that GPT-4 has 1.76 trillion parameters, which was first estimated by the speed it was running and by [George Hotz](#).^[12]

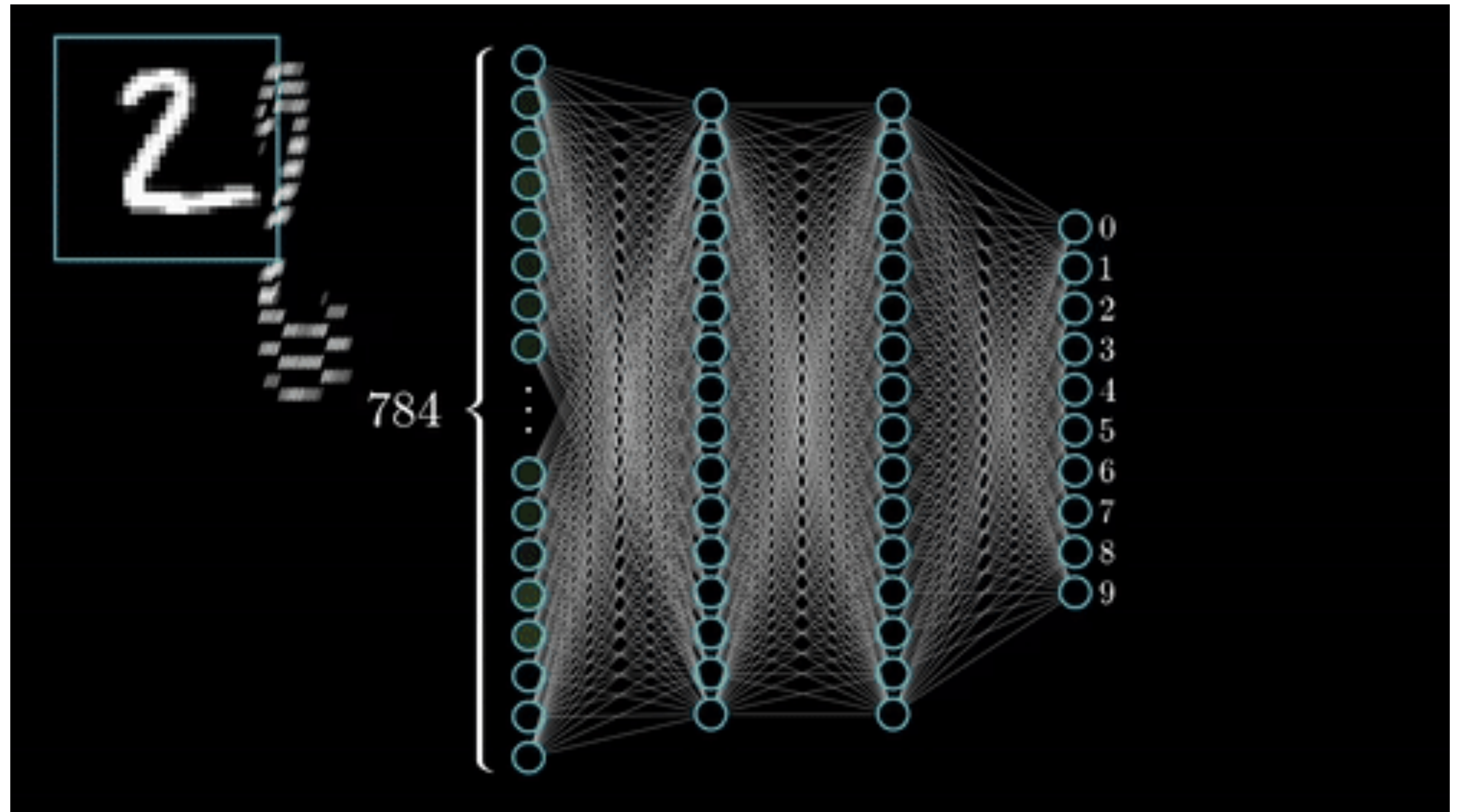
Neural Networks: Backbone of modern ML

Artificial Neural Networks are **highly expressive**

But they have
Large number of parameters

Each edge is a parameter.

Parameters are learned by an algorithm called **Backpropagation**.



Rumors claim that **GPT-4 has 1.76 trillion parameters**, which was first estimated by the speed it was running and by [George Hotz](#).^[12]

Neural Networks: Backbone of modern ML

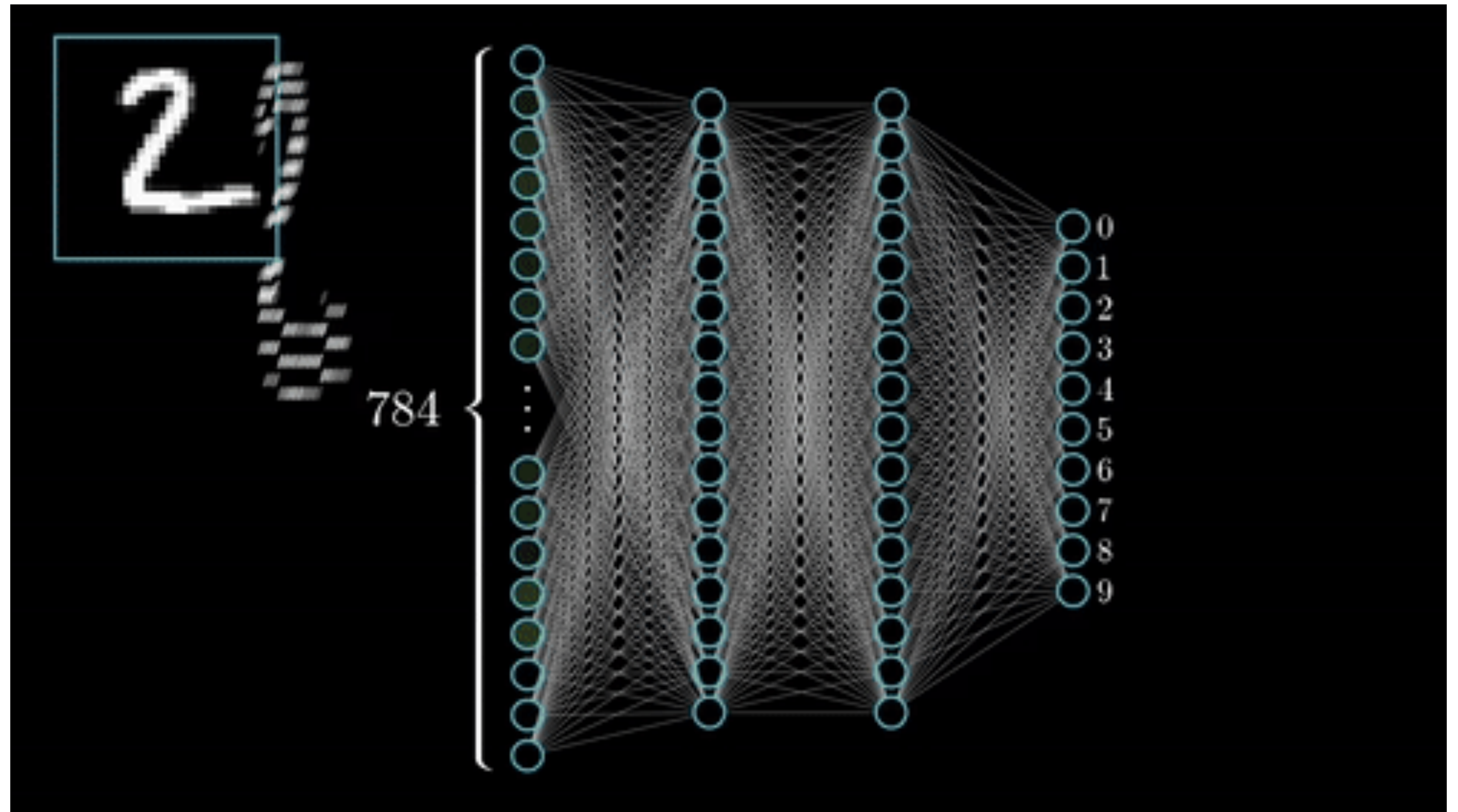
Artificial Neural Networks are **highly expressive**

But they have **Large number of parameters**

Each edge is a parameter.

Parameters are learned by an algorithm called **Backpropagation**.

Backpropagation propagates the discrepancy between the current output and the desired output to all parameters.



Rumors claim that GPT-4 has 1.76 trillion parameters, which was first estimated by the speed it was running and by [George Hotz](#).^[12]

Neural Networks: Backbone of modern ML

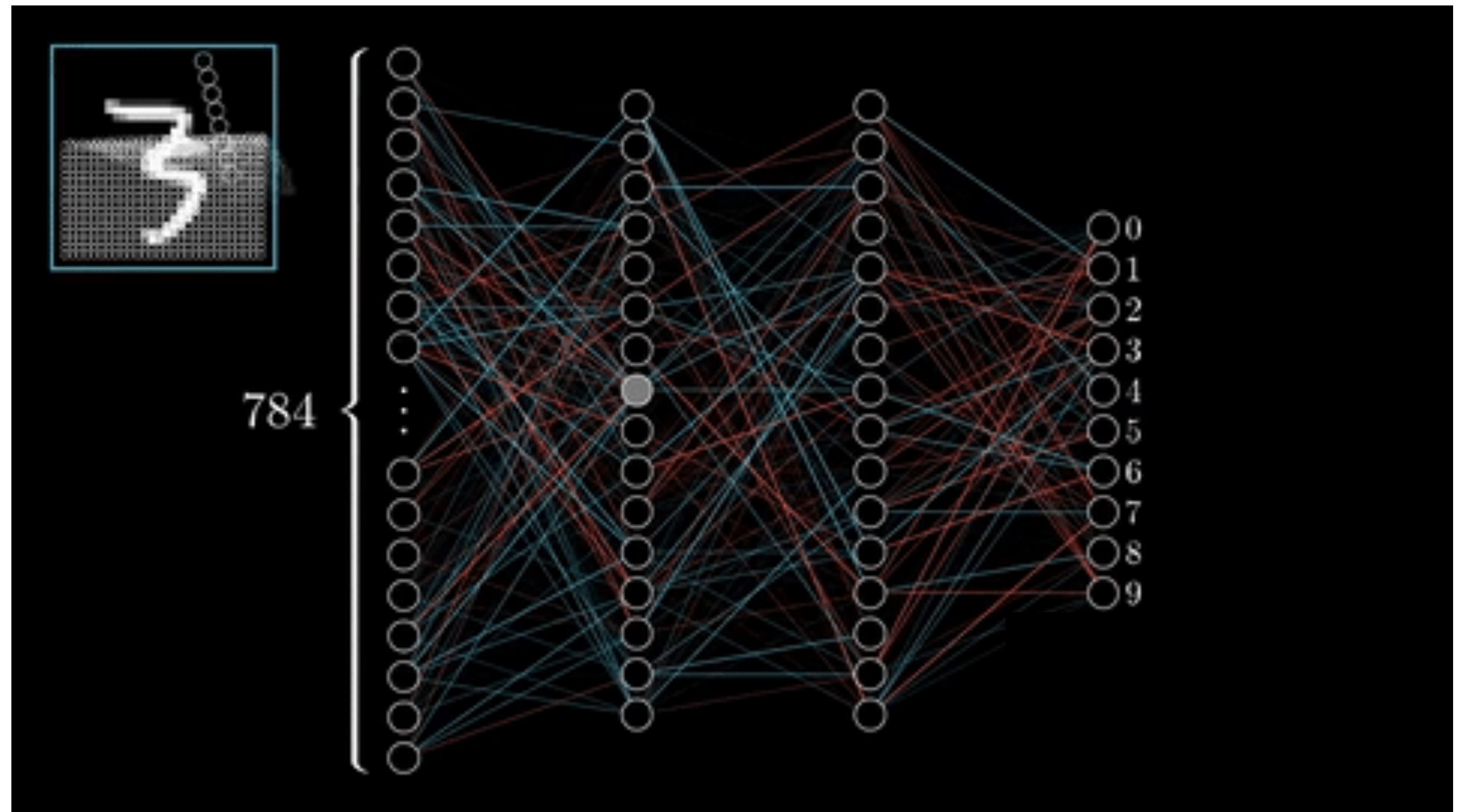
Artificial Neural Networks are **highly expressive**

But they have **Large number of parameters**

Each edge is a parameter.

Parameters are learned by an algorithm called **Backpropagation**.

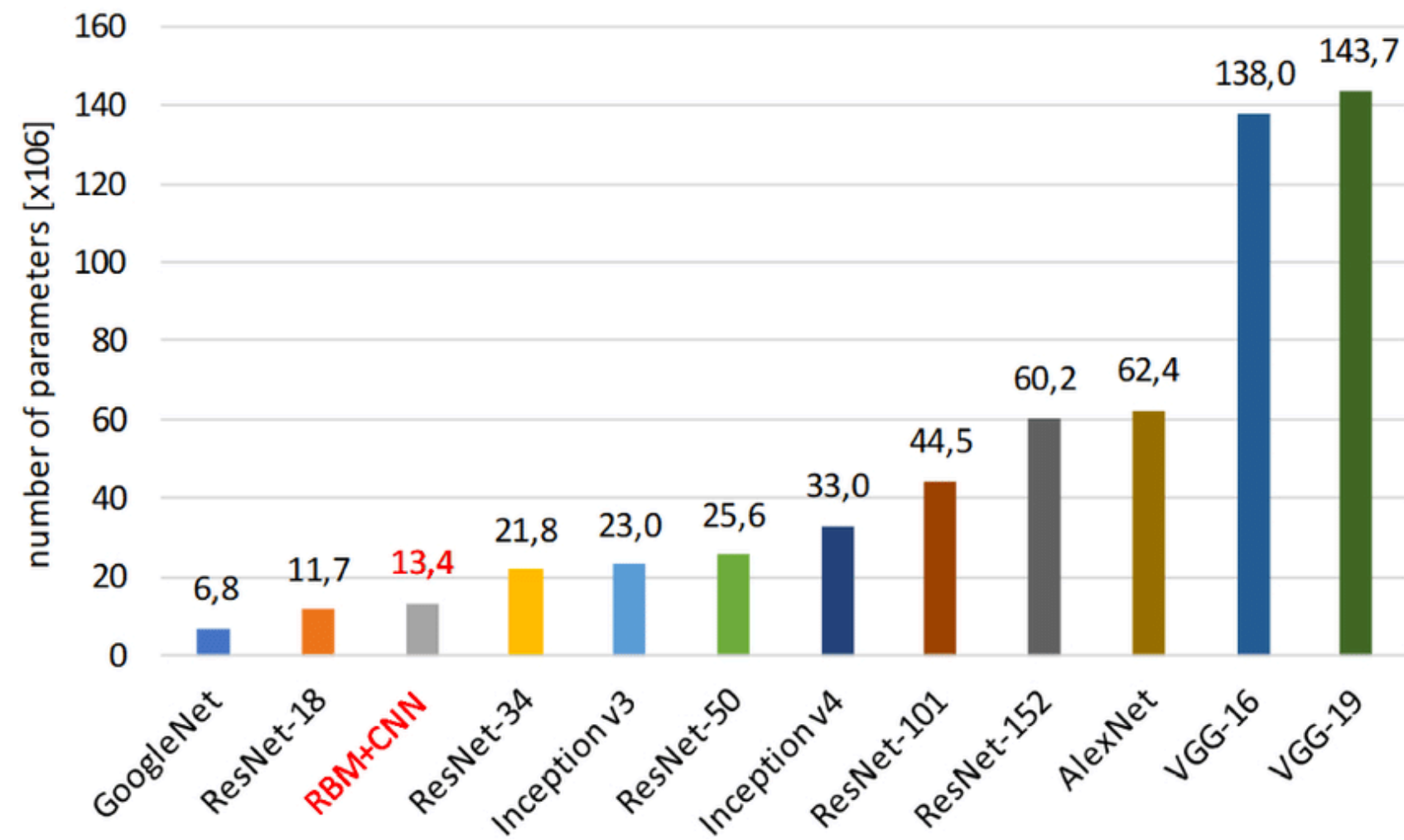
Backpropagation propagates the discrepancy between the current output and the desired output to all parameters.



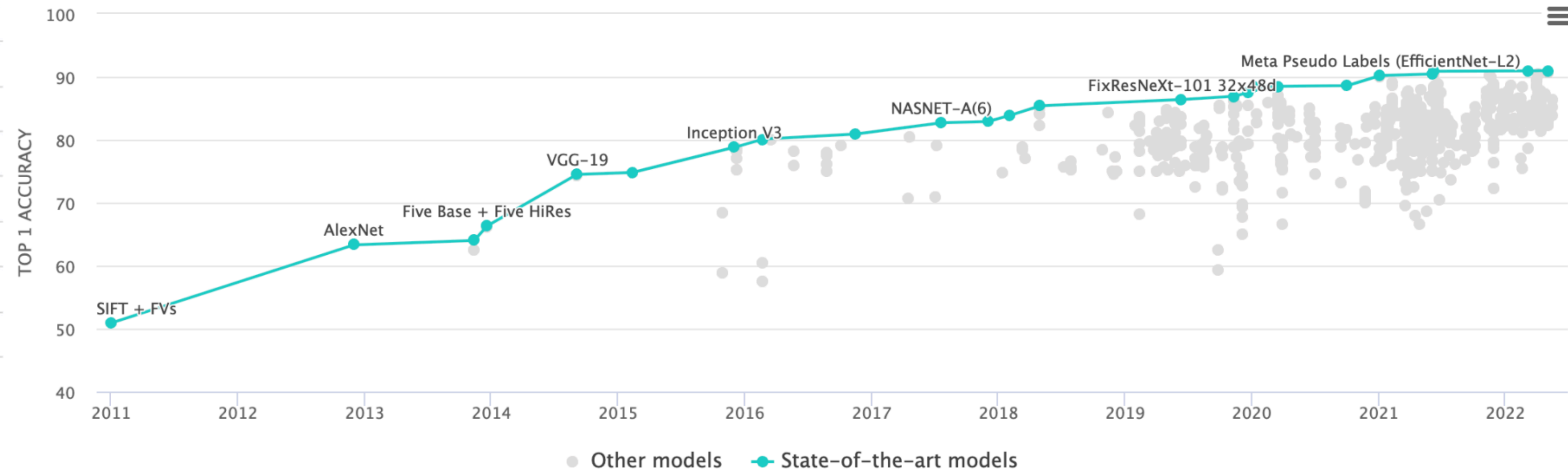
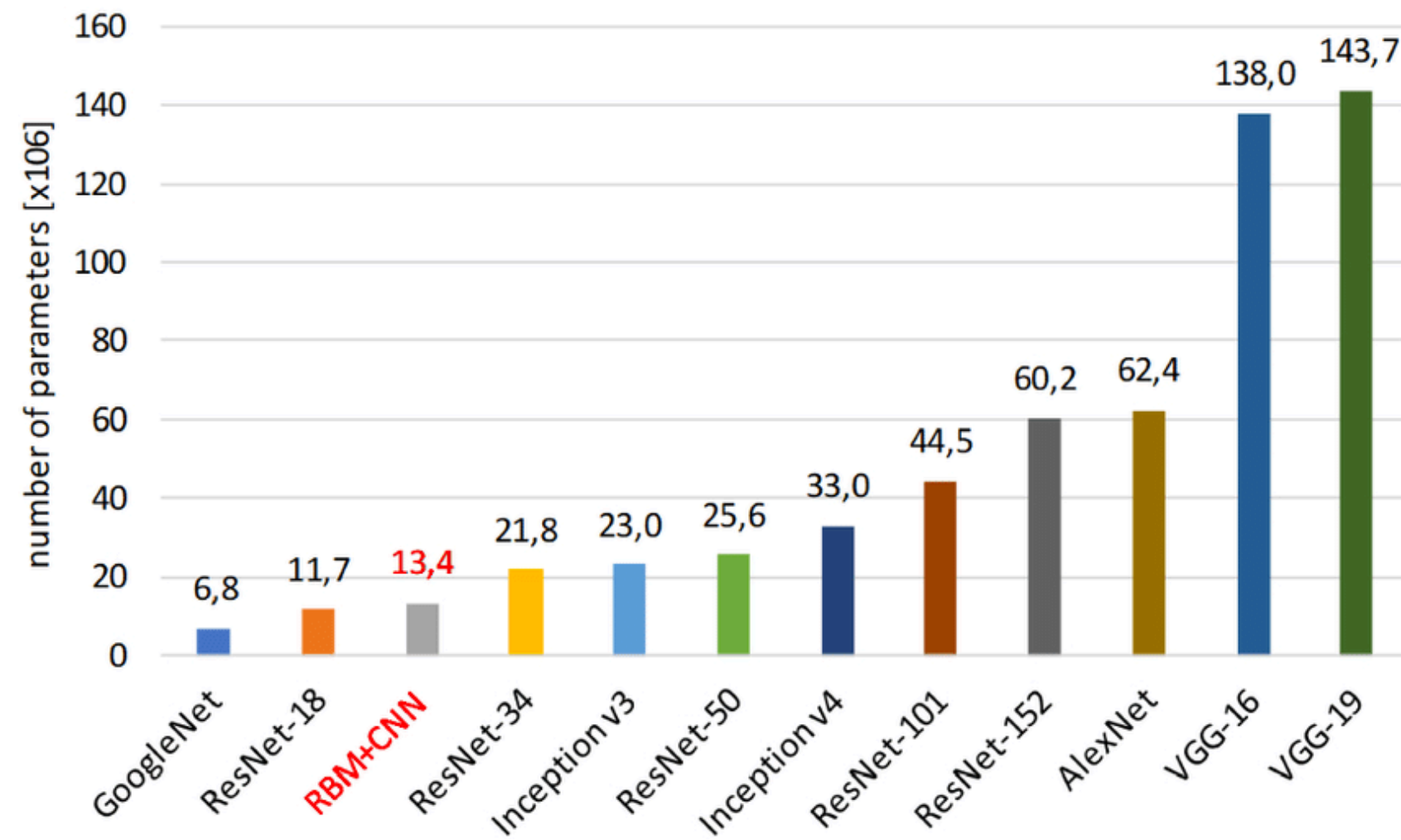
Rumors claim that GPT-4 has 1.76 trillion parameters, which was first estimated by the speed it was running and by [George Hotz](#).^[12]

Growth of Computer Vision models

Growth of Computer Vision models



Growth of Computer Vision models



Application: Language Modelling

Application: Language Modelling

In Computer vision, the task is to **predict the label** given a image.

Application: Language Modelling

In Computer vision, the task is to **predict the label** given a image.

Given a dataset of **images and corresponding labels**, ML algorithm maximises the probability of **predicting the correct label**.

Application: Language Modelling

In Computer vision, the task is to **predict the label** given a image.

Given a dataset of **images and corresponding labels**, ML algorithm maximises the probability of **predicting the correct label**.

In Language modelling, task is to **predict the next token given previous tokens**.
A token can be a character, a word or something in between.

Application: Language Modelling

In Computer vision, the task is to **predict the label** given a image.

Given a dataset of **images and corresponding labels**, ML algorithm maximises the probability of **predicting the correct label**.

In Language modelling, task is to **predict the next token given previous tokens**.
A token can be a character, a word or something in between.

Types of Language Modelling:

Application: Language Modelling

In Computer vision, the task is to **predict the label** given a image.

Given a dataset of **images and corresponding labels**, ML algorithm maximises the probability of **predicting the correct label**.

In Language modelling, task is to **predict the next token given previous tokens**.
A token can be a character, a word or something in between.

Types of Language Modelling:

- **Statistical Language Modelling**

Application: Language Modelling

In Computer vision, the task is to **predict the label** given a image.

Given a dataset of **images and corresponding labels**, ML algorithm maximises the probability of **predicting the correct label**.

In Language modelling, task is to **predict the next token given previous tokens**.
A token can be a character, a word or something in between.

Types of Language Modelling:

- **Statistical Language Modelling**
- **Neural Language Modelling**

Statistical Language Modelling

Statistical Language Modelling

4-gram modelling (generally n-gram)

Statistical Language Modelling

4-gram modelling (generally n-gram)

From a dataset of documents

- **Compute the joint probability** of *all* token sequences (length 4 and 3).

Statistical Language Modelling

4-gram modelling (generally n-gram)

From a dataset of documents

- **Compute the joint probability** of *all* token sequences (length 4 and 3).

We were all feeling seedy, and we were getting quite nervous about it. Harris said he felt such extraordinary fits of giddiness come over him at times, that he hardly knew what he was doing; and then George said that *he* had fits of giddiness too, and hardly knew what *he* was doing. With me, it was my liver that was out of order. I knew it was my liver that was out of order, because I had just been reading a patent liver-pill circular, in which were detailed the various symptoms by which a man could tell when his liver was out of order. I had them all.

Statistical Language Modelling

4-gram modelling (generally n-gram)

From a dataset of documents

- **Compute the joint probability** of *all* token sequences (length 4 and 3).

4-gram	Frequency
was out of order	3
my liver that was	2
it was my liver	2
that was out of	2
liver that was out	2
was my liver that	2

We were all feeling seedy, and we were getting quite nervous about it. Harris said he felt such extraordinary fits of giddiness come over him at times, that he hardly knew what he was doing; and then George said that *he* had fits of giddiness too, and hardly knew what *he* was doing. With me, it was my liver that was out of order. I knew it was my liver that was out of order, because I had just been reading a patent liver-pill circular, in which were detailed the various symptoms by which a man could tell when his liver was out of order. I had them all.

Statistical Language Modelling

4-gram modelling (generally n-gram)

From a dataset of documents

- **Compute the joint probability** of *all* token sequences (length 4 and 3).
- User inputs the first 3 words w_1, w_2, w_3

4-gram	Frequency
was out of order	3
my liver that was	2
it was my liver	2
that was out of	2
liver that was out	2
was my liver that	2

We were all feeling seedy, and we were getting quite nervous about it. Harris said he felt such extraordinary fits of giddiness come over him at times, that he hardly knew what he was doing; and then George said that *he* had fits of giddiness too, and hardly knew what *he* was doing. With me, it was my liver that was out of order. I knew it was my liver that was out of order, because I had just been reading a patent liver-pill circular, in which were detailed the various symptoms by which a man could tell when his liver was out of order. I had them all.

Statistical Language Modelling

4-gram modelling (generally n-gram)

From a dataset of documents

- **Compute the joint probability** of *all* token sequences (length 4 and 3).
- User inputs the first 3 words w_1, w_2, w_3
- Then given w_1, w_2, w_3 , **sample** w_4 **from the conditional distribution** $p(w \mid w_1, w_2, w_3)$.

4-gram	Frequency
was out of order	3
my liver that was	2
it was my liver	2
that was out of	2
liver that was out	2
was my liver that	2

We were all feeling seedy, and we were getting quite nervous about it. Harris said he felt such extraordinary fits of giddiness come over him at times, that he hardly knew what he was doing; and then George said that *he* had fits of giddiness too, and hardly knew what *he* was doing. With me, it was my liver that was out of order. I knew it was my liver that was out of order, because I had just been reading a patent liver-pill circular, in which were detailed the various symptoms by which a man could tell when his liver was out of order. I had them all.

Statistical Language Modelling

4-gram modelling (generally n-gram)

From a dataset of documents

- **Compute the joint probability** of *all* token sequences (length 4 and 3).
- User inputs the first 3 words w_1, w_2, w_3
- Then given w_1, w_2, w_3 , **sample** w_4 **from the conditional distribution** $p(w \mid w_1, w_2, w_3)$.
- Then keep repeating the process.

4-gram	Frequency
was out of order	3
my liver that was	2
it was my liver	2
that was out of	2
liver that was out	2
was my liver that	2

We were all feeling seedy, and we were getting quite nervous about it. Harris said he felt such extraordinary fits of giddiness come over him at times, that he hardly knew what he was doing; and then George said that *he* had fits of giddiness too, and hardly knew what *he* was doing. With me, it was my liver that was out of order. I knew it was my liver that was out of order, because I had just been reading a patent liver-pill circular, in which were detailed the various symptoms by which a man could tell when his liver was out of order. I had them all.

Statistical Language Modelling

4-gram modelling (generally n-gram)

From a dataset of documents

- **Compute the joint probability** of *all* token sequences (length 4 and 3).
- User inputs the first 3 words w_1, w_2, w_3
- Then given w_1, w_2, w_3 , **sample** w_4 **from the conditional distribution** $p(w \mid w_1, w_2, w_3)$.
- Then keep repeating the process.

4-gram	Frequency
was out of order	3
my liver that was	2
it was my liver	2
that was out of	2
liver that was out	2
was my liver that	2

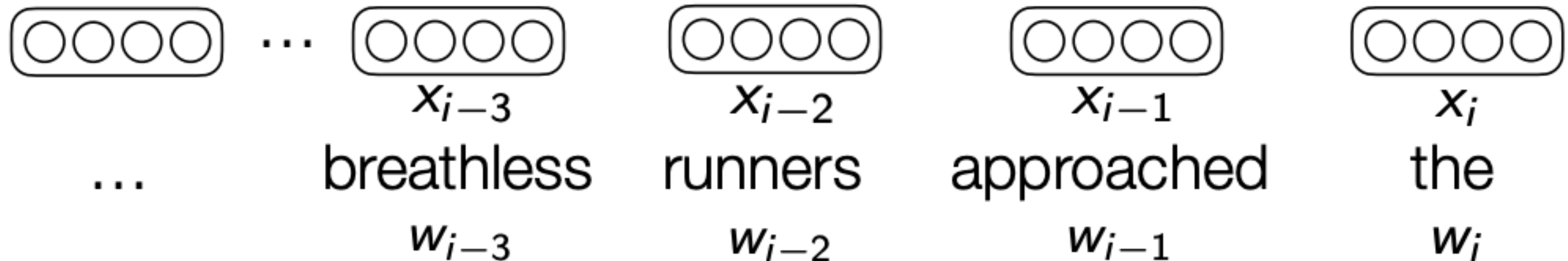
We were all feeling seedy, and we were getting quite nervous about it. Harris said he felt such extraordinary fits of giddiness come over him at times, that he hardly knew what he was doing; and then George said that *he* had fits of giddiness too, and hardly knew what *he* was doing. With me, it was my liver that was out of order. I knew it was my liver that was out of order, because I had just been reading a patent liver-pill circular, in which were detailed the various symptoms by which a man could tell when his liver was out of order. I had them all.

This is a simple process but it cannot model long-term dependencies.

Neural Language Modelling

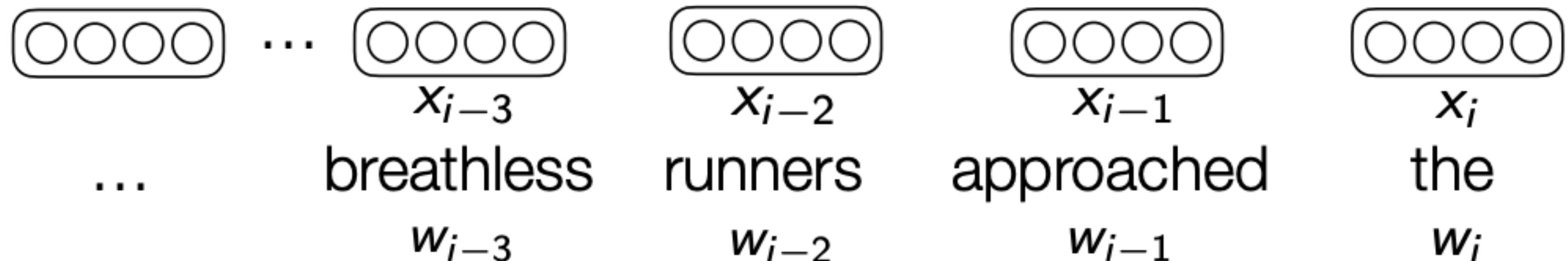
Neural Language Modelling

- Input: Words w_i



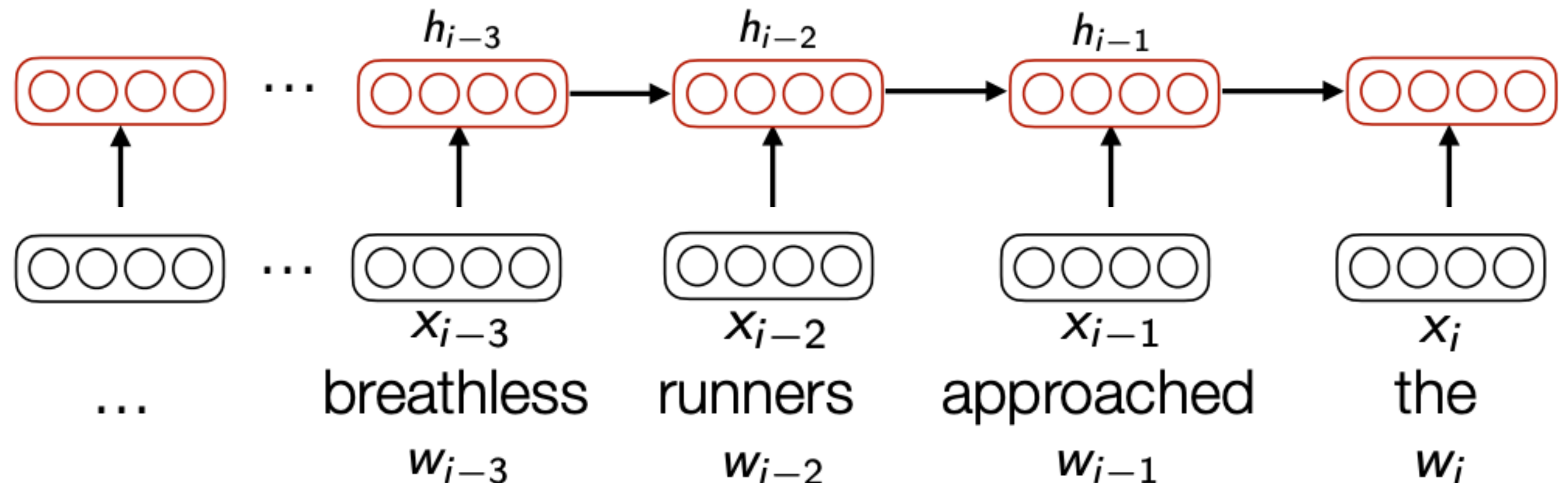
Neural Language Modelling

- Input: Words w_i
- At each time step, the neural network does two things



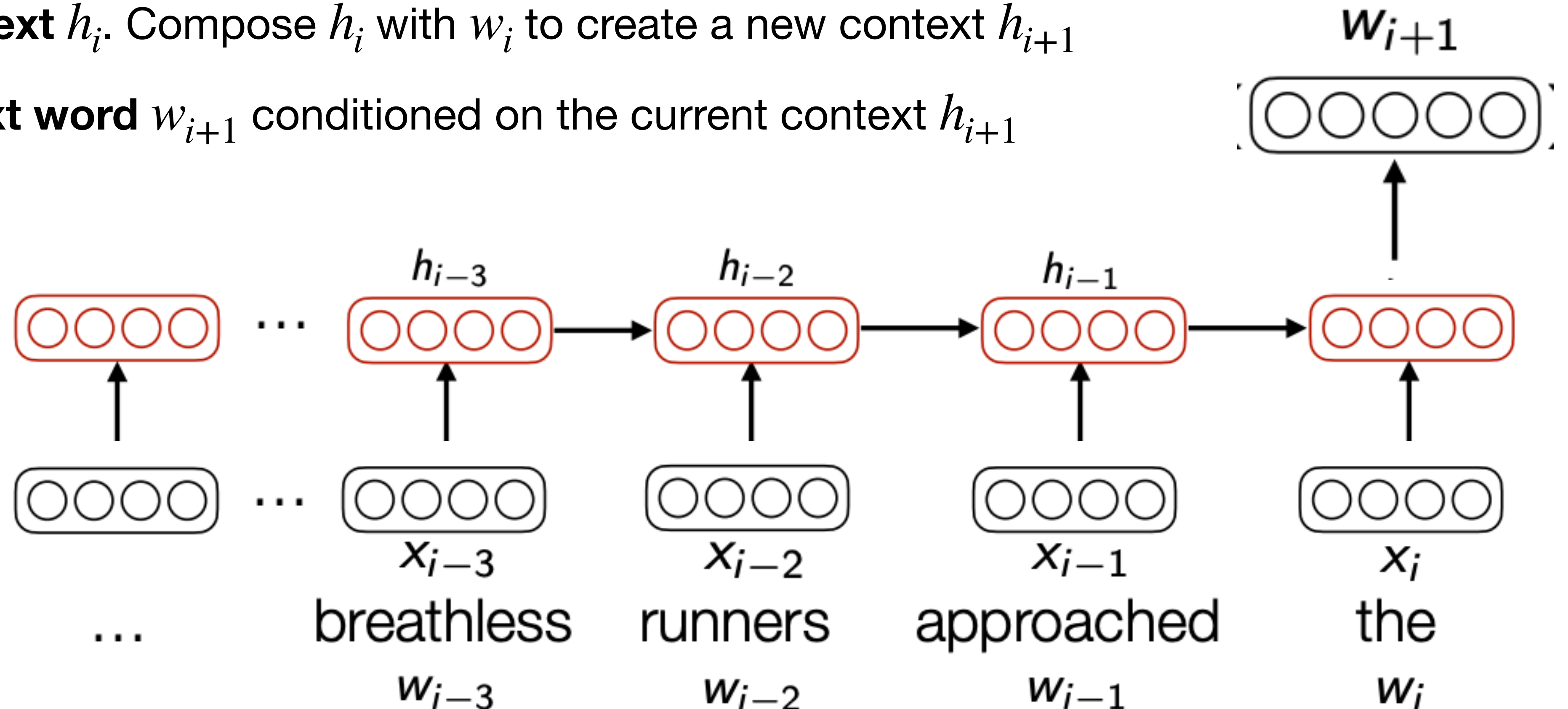
Neural Language Modelling

- Input: Words w_i
- At each time step, the neural network does two things
 - **Maintain a context** h_i . Compose h_i with w_i to create a new context h_{i+1}



Neural Language Modelling

- Input: Words w_i
- At each time step, the neural network does two things
 - **Maintain a context** h_i . Compose h_i with w_i to create a new context h_{i+1}
 - **Predicts the next word** w_{i+1} conditioned on the current context h_{i+1}



Modern Machine Learning: The good

Modern Machine Learning: The good

We are able to train large Machine Learning models that surpass average humans in various tasks.

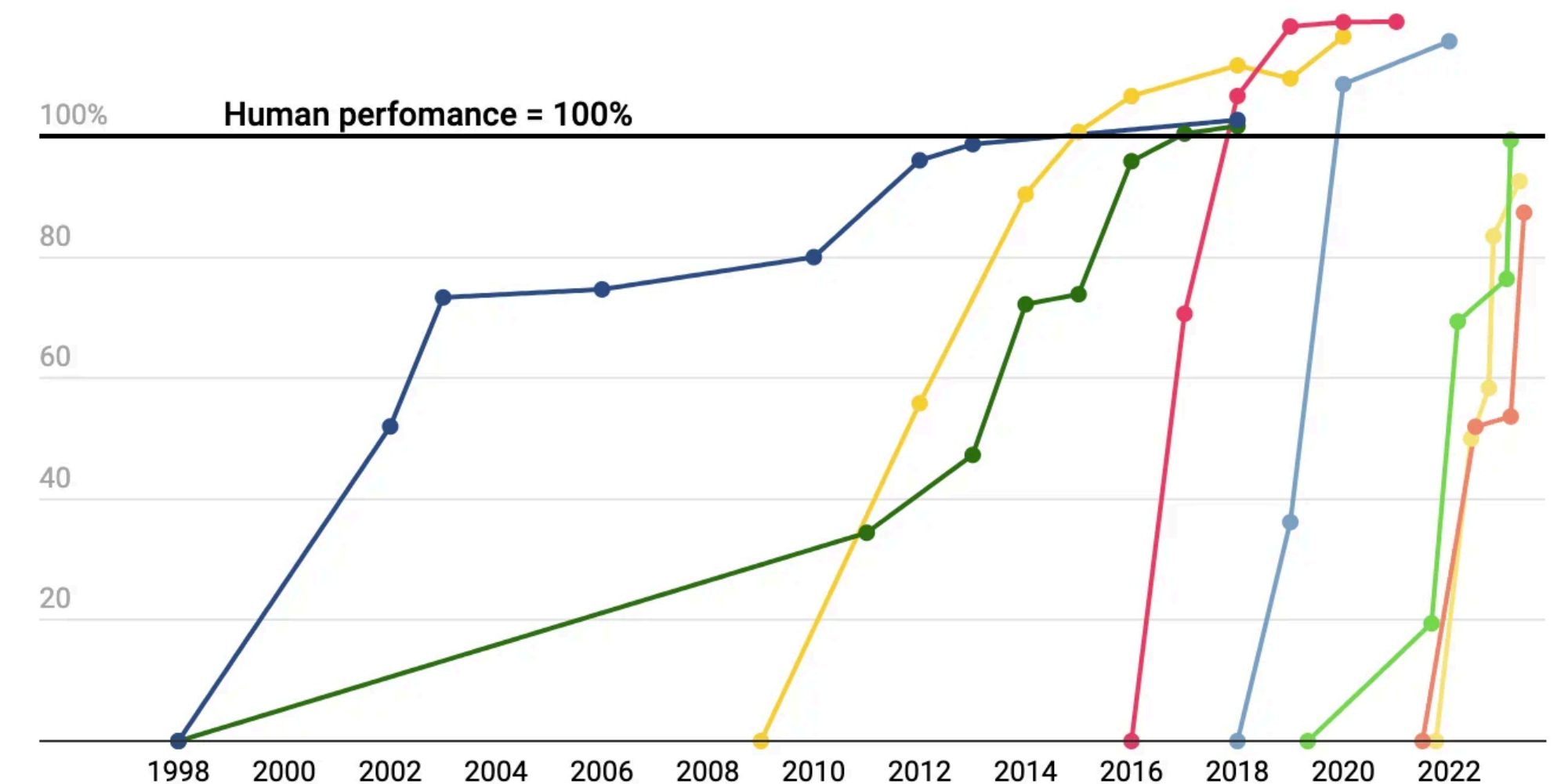
Modern Machine Learning: The good

We are able to train large Machine Learning models that surpass average humans in various tasks.

AI has surpassed humans at a number of tasks and the rate at which humans are being surpassed at new tasks is increasing

State-of-the-art AI performance on benchmarks, relative to human performance

● Handwriting recognition ● Speech recognition ● Image recognition ● Reading comprehension
● Language understanding ● Common sense completion ● Grade school math ● Code generation



For each benchmark, the maximally performing baseline reported in the benchmark paper is taken as the "starting point", which is set at 0%. Human performance number is set at 100%. Handwriting recognition = MNIST, Language understanding = GLUE, Image recognition = ImageNet, Reading comprehension = SQuAD 1.1, Reading comprehension = SQuAD 2.0, Speech recognition = Switchboard, Grade school math = GSK8k, Common sense completion = HellaSwag, Code generation = HumanEval.

Chart: Will Henshall for TIME • Source: [ContextualAI](#)

TIME

Modern Machine Learning: The good

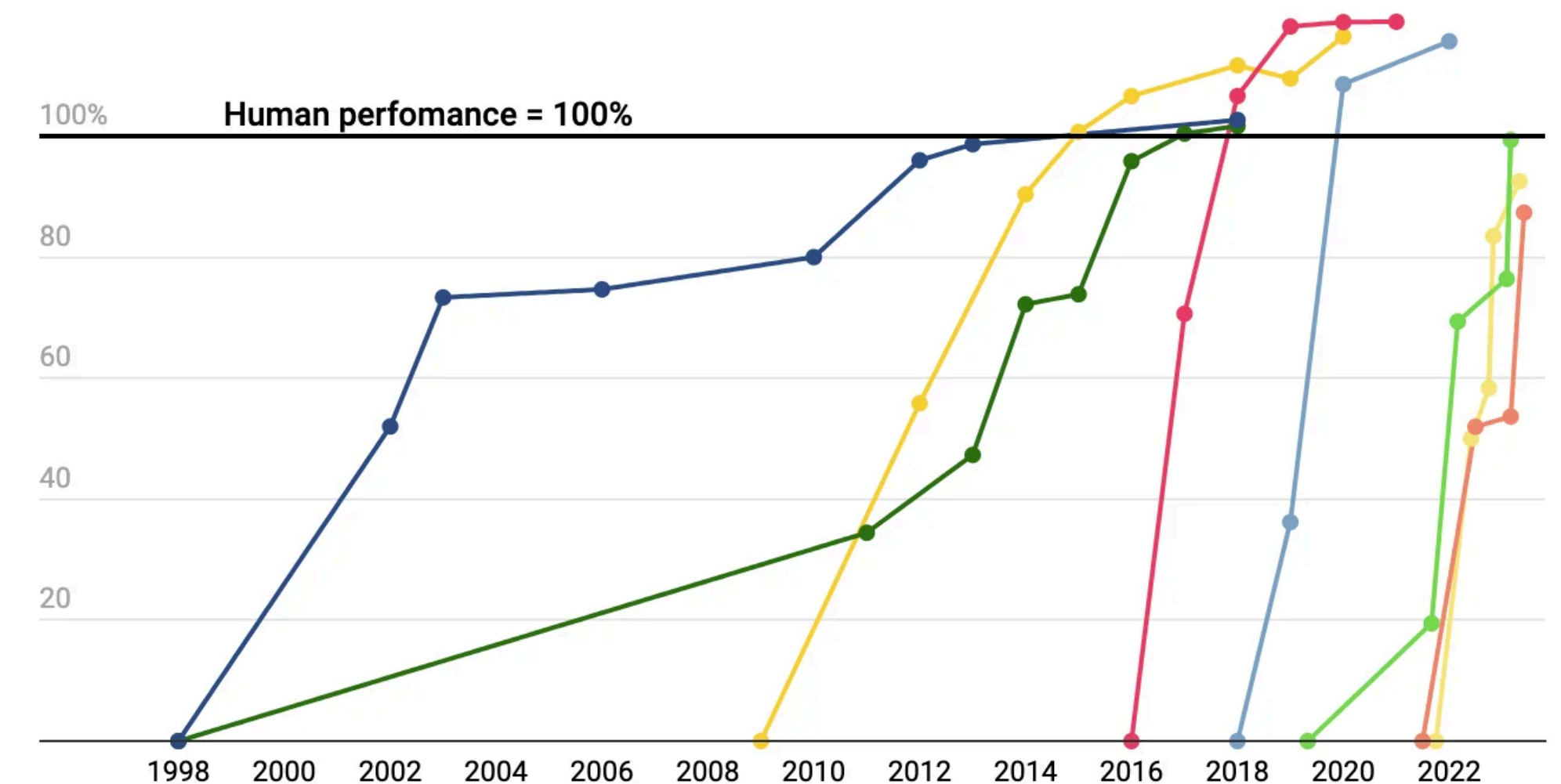
We are able to train large Machine Learning models that surpass average humans in various tasks.



AI has surpassed humans at a number of tasks and the rate at which humans are being surpassed at new tasks is increasing

State-of-the-art AI performance on benchmarks, relative to human performance

● Handwriting recognition ● Speech recognition ● Image recognition ● Reading comprehension
● Language understanding ● Common sense completion ● Grade school math ● Code generation



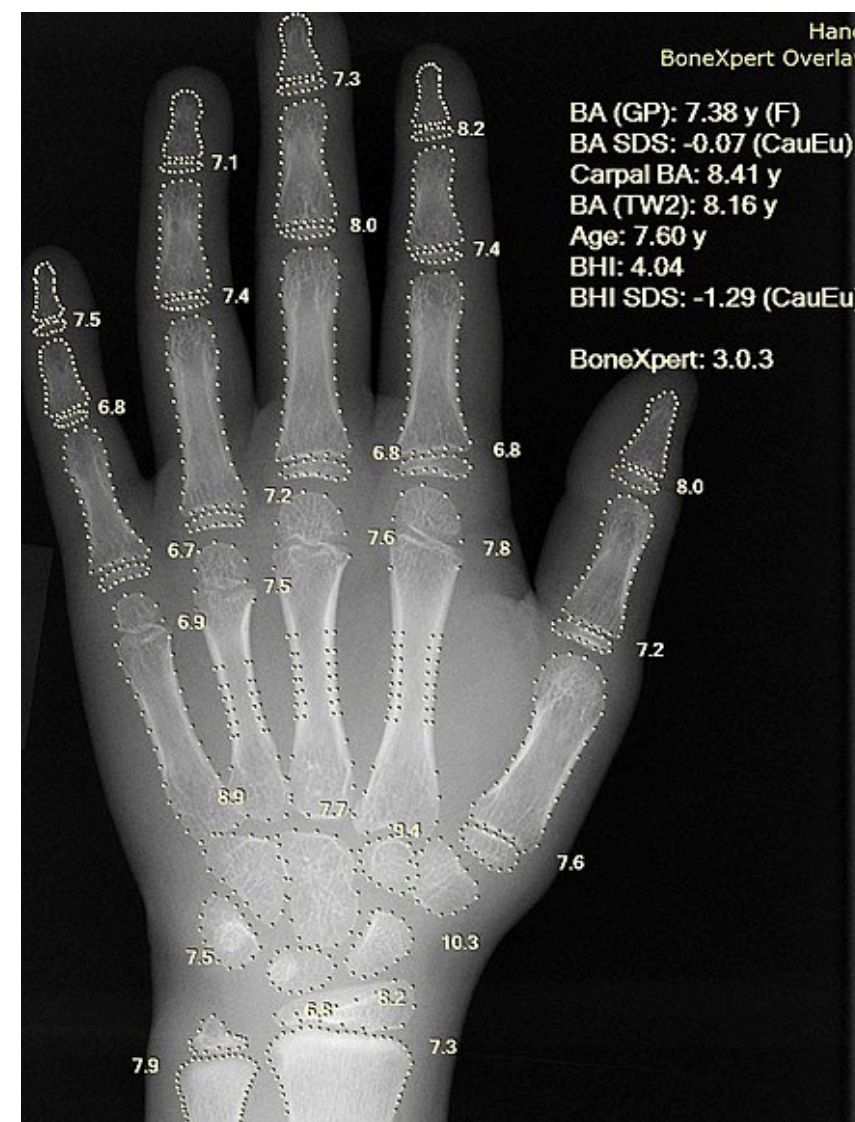
For each benchmark, the maximally performing baseline reported in the benchmark paper is taken as the "starting point", which is set at 0%. Human performance number is set at 100%. Handwriting recognition = MNIST, Language understanding = GLUE, Image recognition = ImageNet, Reading comprehension = SQuAD 1.1, Reading comprehension = SQuAD 2.0, Speech recognition = Switchboard, Grade school math = GSK8k, Common sense completion = HellaSwag, Code generation = HumanEval.

Chart: Will Henshall for TIME • Source: [ContextualAI](#)

TIME

Modern Machine Learning: The good

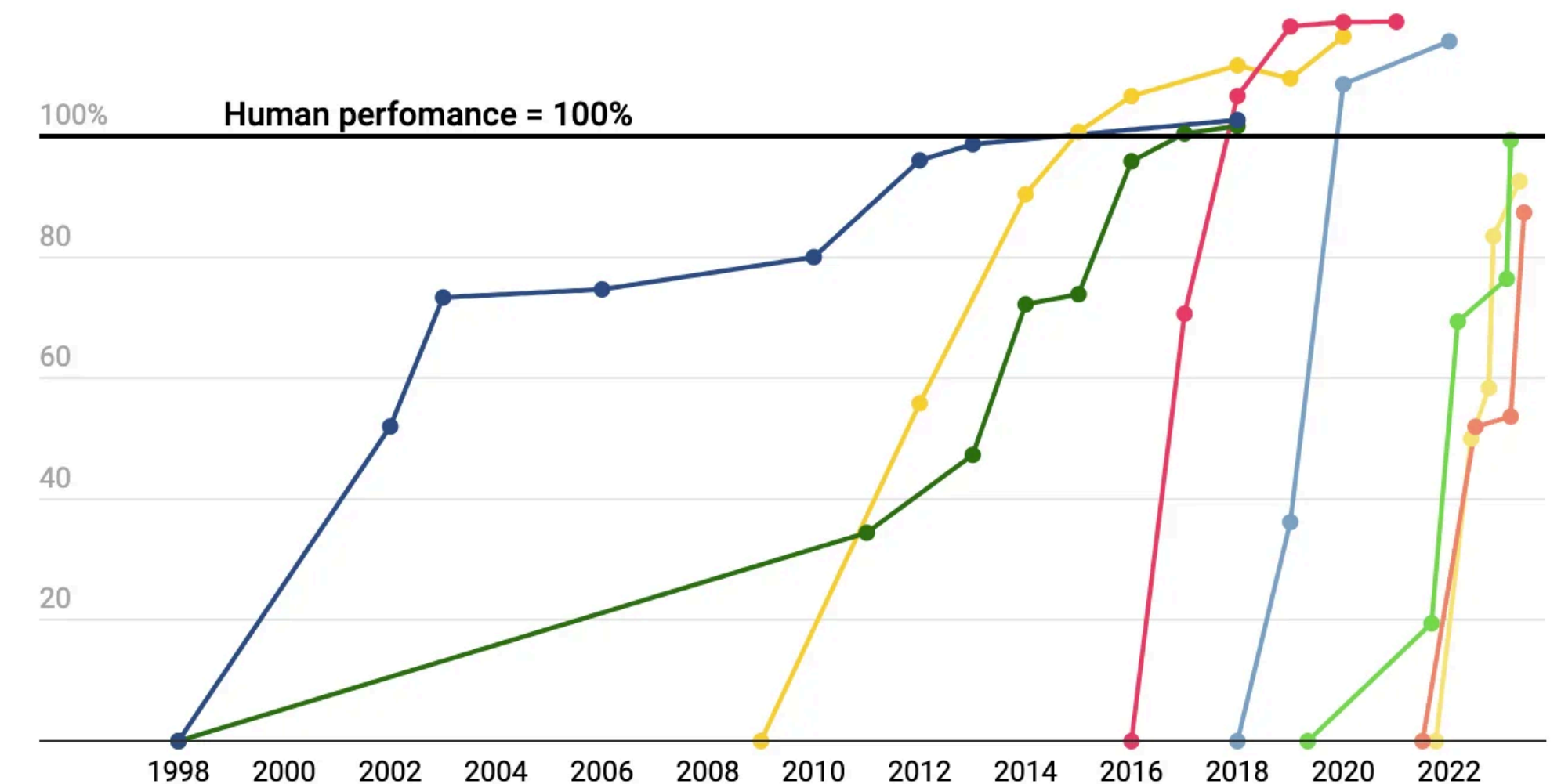
We are able to train large Machine Learning models that surpass average humans in various tasks.



AI has surpassed humans at a number of tasks and the rate at which humans are being surpassed at new tasks is increasing

State-of-the-art AI performance on benchmarks, relative to human performance

● Handwriting recognition ● Speech recognition ● Image recognition ● Reading comprehension
● Language understanding ● Common sense completion ● Grade school math ● Code generation



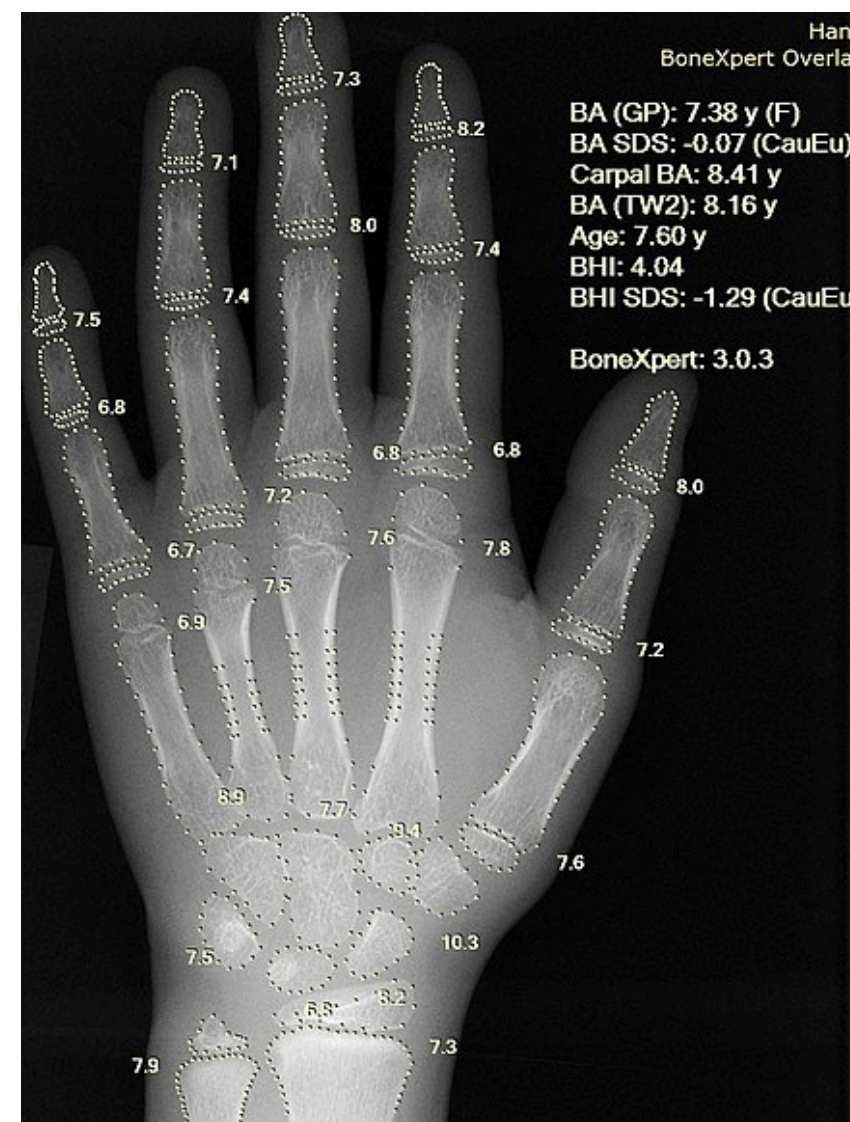
For each benchmark, the maximally performing baseline reported in the benchmark paper is taken as the "starting point", which is set at 0%. Human performance number is set at 100%. Handwriting recognition = MNIST, Language understanding = GLUE, Image recognition = ImageNet, Reading comprehension = SQuAD 1.1, Reading comprehension = SQuAD 2.0, Speech recognition = Switchboard, Grade school math = GSK8k, Common sense completion = HellaSwag, Code generation = HumanEval.

Chart: Will Henshall for TIME • Source: [ContextualAI](#)

TIME

Modern Machine Learning: The good

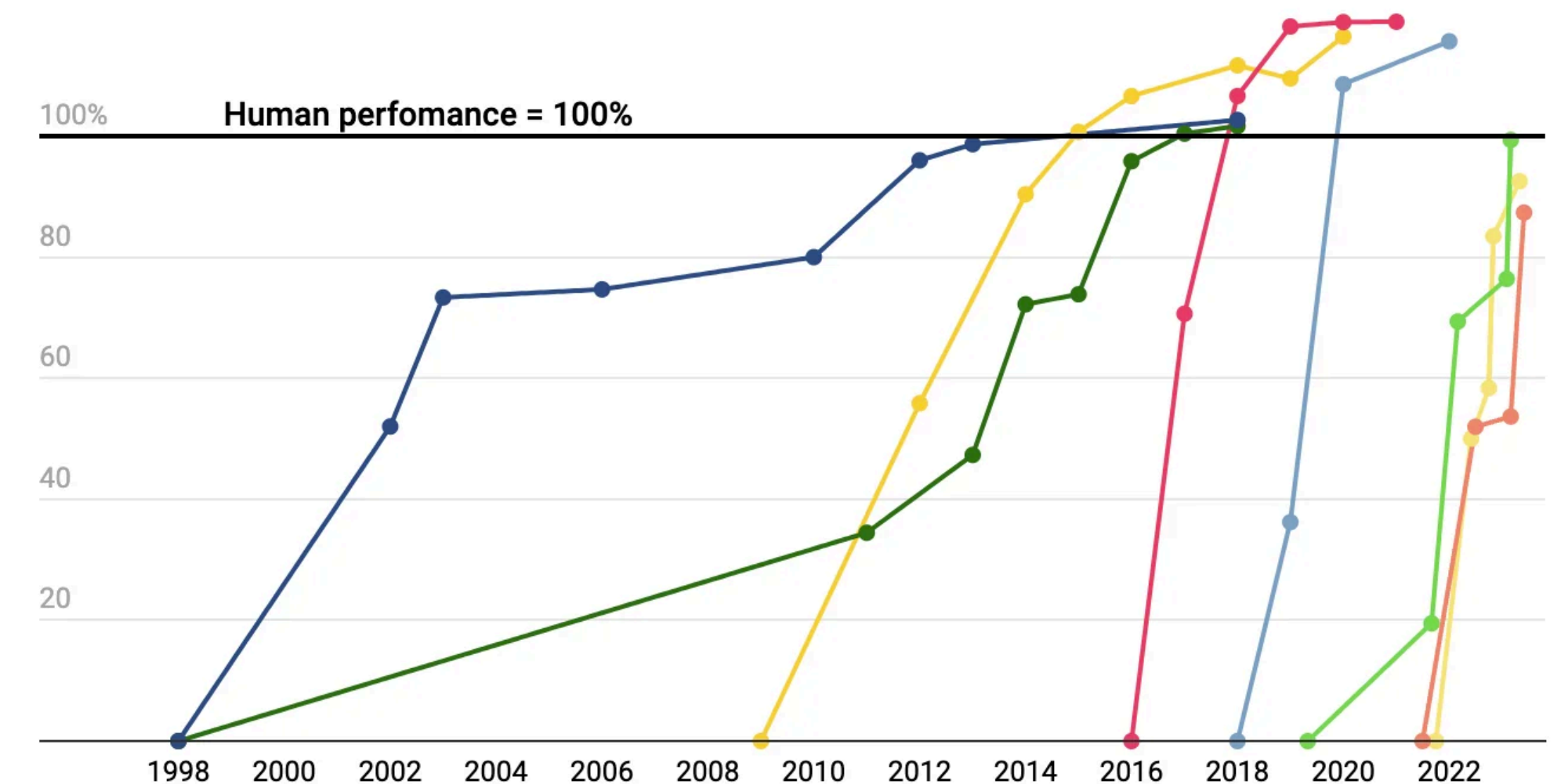
We are able to train large Machine Learning models that surpass average humans in various tasks.



AI has surpassed humans at a number of tasks and the rate at which humans are being surpassed at new tasks is increasing

State-of-the-art AI performance on benchmarks, relative to human performance

● Handwriting recognition ● Speech recognition ● Image recognition ● Reading comprehension
● Language understanding ● Common sense completion ● Grade school math ● Code generation



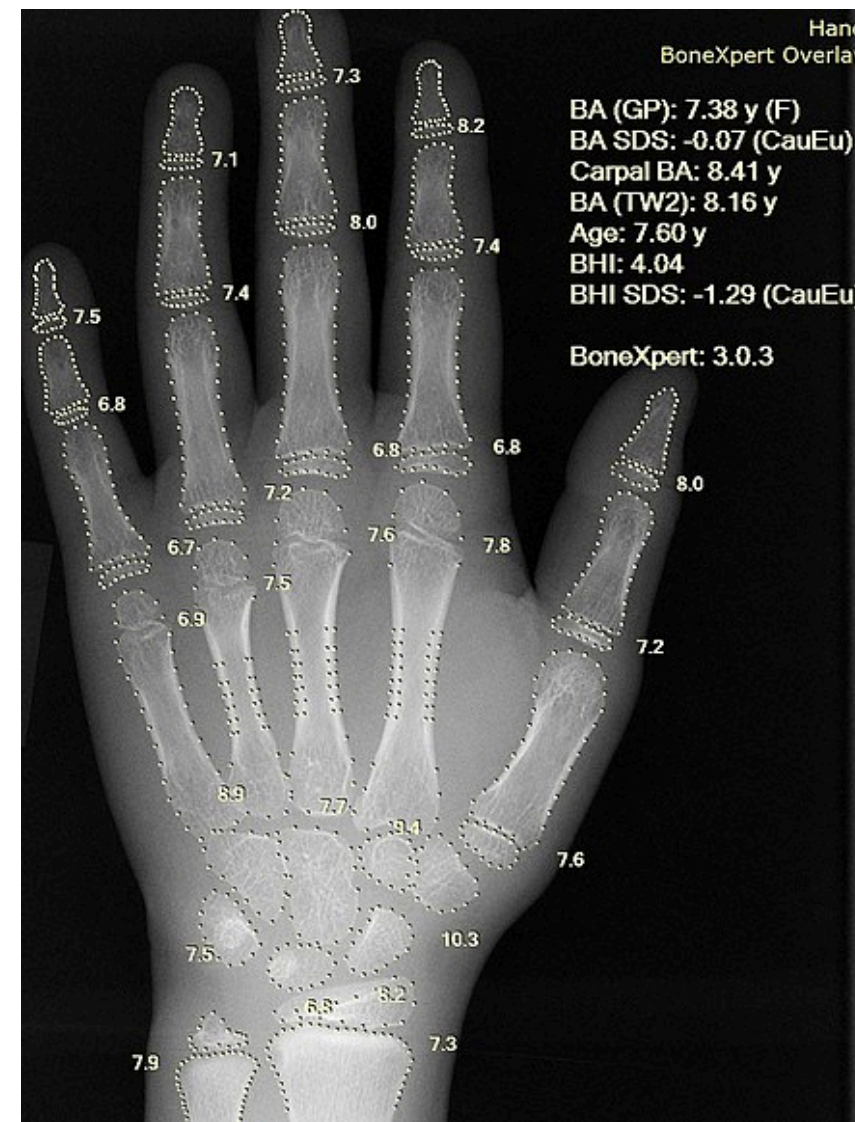
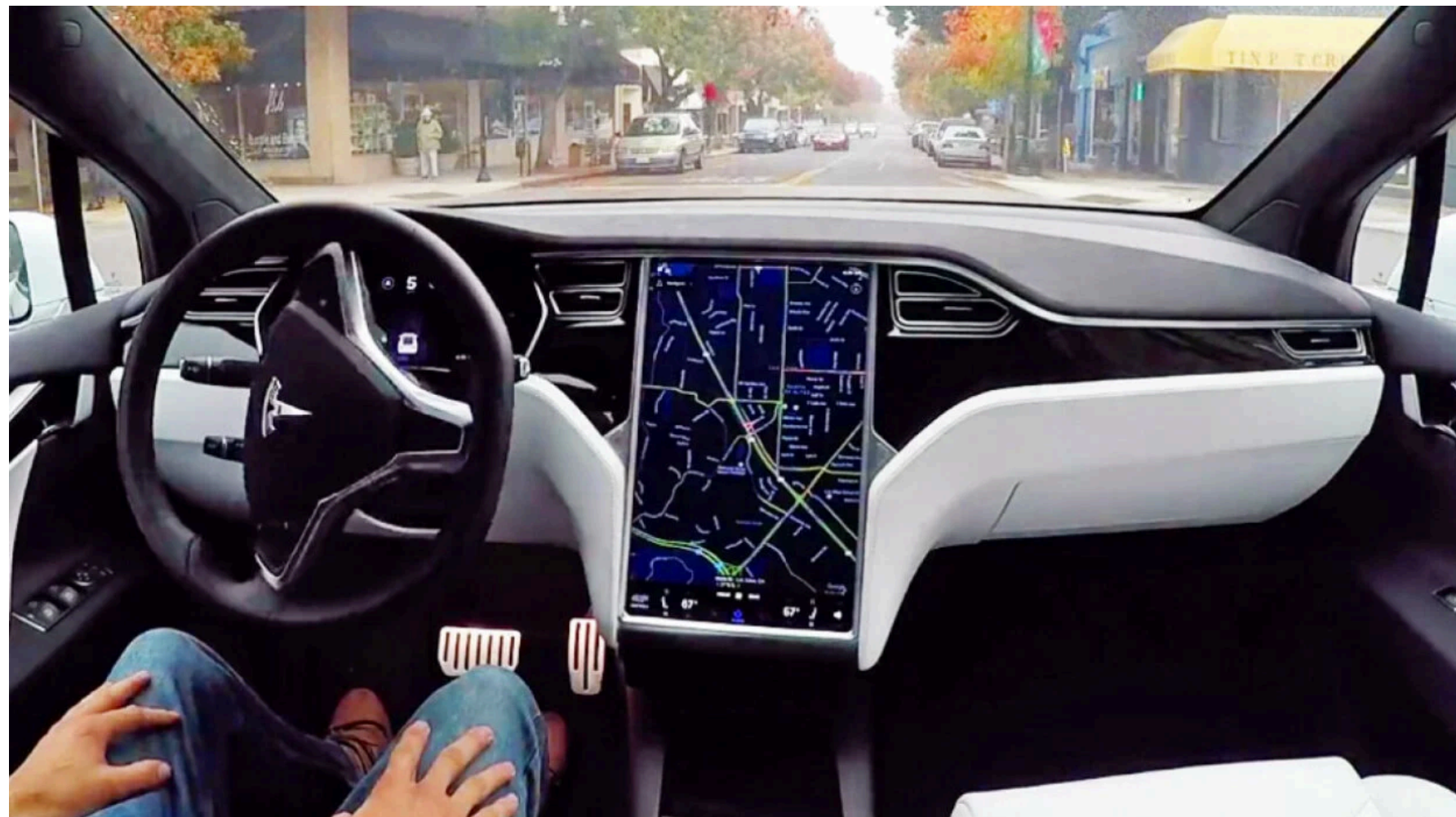
For each benchmark, the maximally performing baseline reported in the benchmark paper is taken as the "starting point", which is set at 0%. Human performance number is set at 100%. Handwriting recognition = MNIST, Language understanding = GLUE, Image recognition = ImageNet, Reading comprehension = SQuAD 1.1, Reading comprehension = SQuAD 2.0, Speech recognition = Switchboard, Grade school math = GSK8k, Common sense completion = HellaSwag, Code generation = HumanEval.

Chart: Will Henshall for TIME • Source: [ContextualAI](#)

TIME

Modern Machine Learning: The good

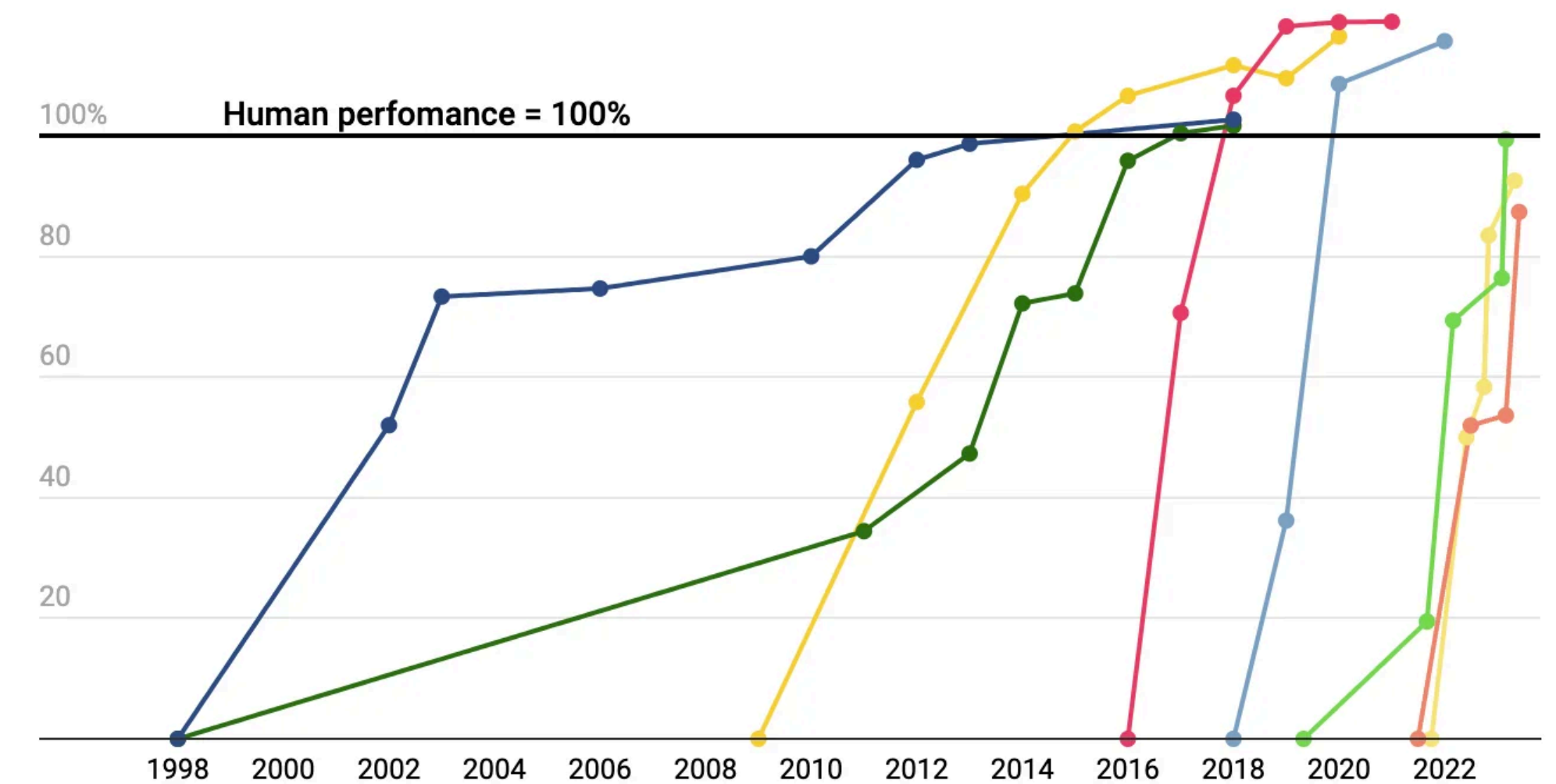
We are able to train large Machine Learning models that surpass average humans in various tasks.



AI has surpassed humans at a number of tasks and the rate at which humans are being surpassed at new tasks is increasing

State-of-the-art AI performance on benchmarks, relative to human performance

● Handwriting recognition ● Speech recognition ● Image recognition ● Reading comprehension
● Language understanding ● Common sense completion ● Grade school math ● Code generation



For each benchmark, the maximally performing baseline reported in the benchmark paper is taken as the "starting point", which is set at 0%. Human performance number is set at 100%. Handwriting recognition = MNIST, Language understanding = GLUE, Image recognition = ImageNet, Reading comprehension = SQuAD 1.1, Reading comprehension = SQuAD 2.0, Speech recognition = Switchboard, Grade school math = GSK8k, Common sense completion = HellaSwag, Code generation = HumanEval.

Chart: Will Henshall for TIME • Source: [ContextualAI](#)

TIME

Modern Machine Learning: The bad

Modern Machine Learning: The bad

To train large machine learning models, we need
Lots of High Quality User relevant training data.

Modern Machine Learning: The bad

To train large machine learning models, we need
Lots of High Quality User relevant training data.

This leads to various concerns in



Privacy

Modern Machine Learning: The bad

To train large machine learning models, we need
Lots of High Quality User relevant training data.

This leads to various concerns in



Privacy



Fairness

Modern Machine Learning: The bad

To train large machine learning models, we need
Lots of High Quality User relevant training data.

This leads to various concerns in

Privacy

Fairness

Robustness

Modern Machine Learning: The bad

To train large machine learning models, we need
Lots of High Quality User relevant training data.

This leads to various concerns in

Privacy

Fairness

Robustness

Backdoors

Image Generation

Image Generation

Machine Learning models can be used to generate images

Image Generation

Machine Learning models can be used to generate images



Image Generation

Machine Learning models can be used to generate images



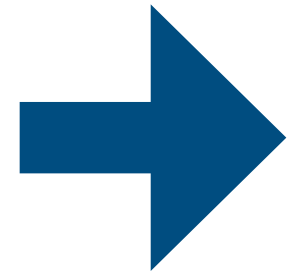
An astronut on a
horse on Mars

Image Generation

Machine Learning models can be used to generate images



An astronut on a
horse on Mars



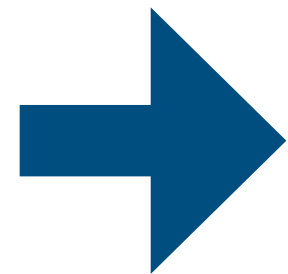
Stable Diffusion

Image Generation

Machine Learning models can be used to generate images



An astronut on a
horse on Mars



Stable Diffusion

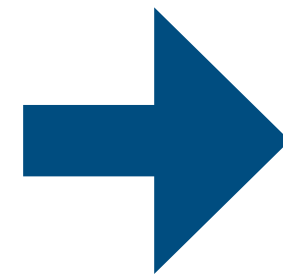
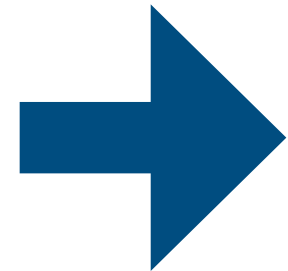


Image Generation

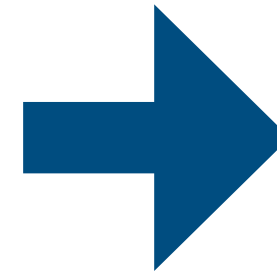
Machine Learning models can be used to generate images



An astronut on a
horse on Mars



Stable Diffusion



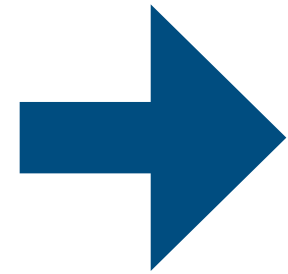
A dog in a rocket aiming
going to the grocery

Image Generation

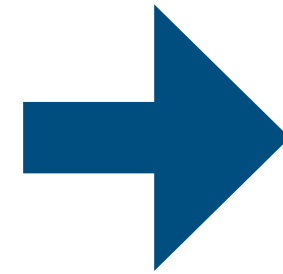
Machine Learning models can be used to generate images



An astronut on a
horse on Mars



Stable Diffusion



A dog in a rocket aiming
going to the grocery



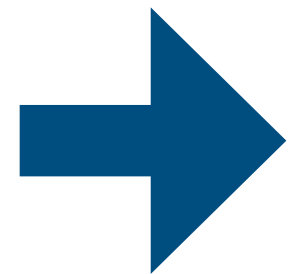
Dall-E

Image Generation

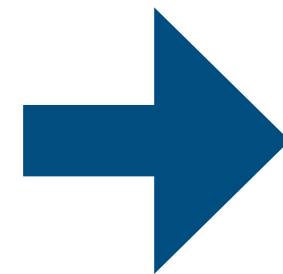
Machine Learning models can be used to generate images



An astronaut on a horse on Mars



Stable Diffusion



A dog in a rocket aiming going to the grocery



Dall-E



Privacy

Privacy

But this data is not always “generated”. They come from training data.

Privacy

But this data is not always “generated”. They come from training data.

Original:



Generated:



Privacy

But this data is not always “generated”. They come from training data.

Has implications for copyright laws.

Original:



Generated:



Text Generation (LLM)

Text Generation (LLM)

Machine Learning models can be used to generate text

Text Generation (LLM)

Machine Learning models can be used to generate text



Me: Write something clever about education.

Text Generation (LLM)

Machine Learning models can be used to generate text



Me: Write something clever about education.

ChatGPT: "Books are the whispering trees of wisdom, where each leaf turned is a step into a forest of knowledge that grows beneath the moonlight of curiosity."

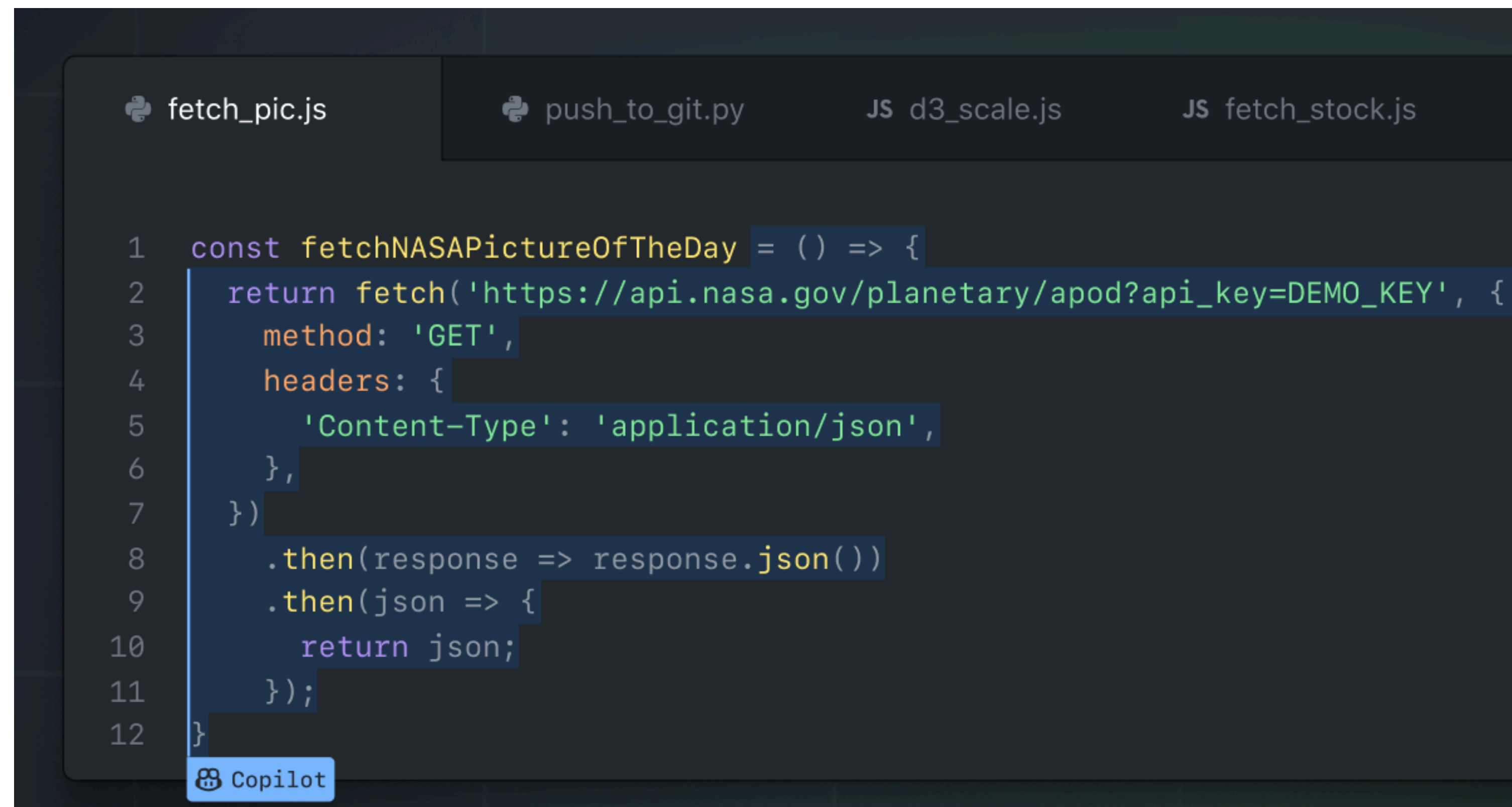
Text Generation (LLM)

Machine Learning models can be used to generate text



Me: Write something clever about education.

ChatGPT: "Books are the whispering trees of wisdom, where each leaf turned is a step into a forest of knowledge that grows beneath the moonlight of curiosity."

A screenshot of the GitHub Copilot interface showing a code editor with a JavaScript file named "fetch_pic.js". The code is a function "fetchNASAPictureOfTheDay" that uses the "fetch" API to retrieve a NASA picture of the day. The code is highlighted in blue. The interface includes a tab bar at the top with several files: "fetch_pic.js", "push_to_git.py", "d3_scale.js", and "fetch_stock.js". The "fetch_pic.js" tab is active. The code is as follows:

```
1  const fetchNASAPictureOfTheDay = () => {
2    return fetch('https://api.nasa.gov/planetary/apod?api_key=DEMO_KEY', {
3      method: 'GET',
4      headers: {
5        'Content-Type': 'application/json',
6      },
7    })
8    .then(response => response.json())
9    .then(json => {
10      return json;
11    });
12  }
```

Privacy

Privacy

Training data can be extracted from large language models.

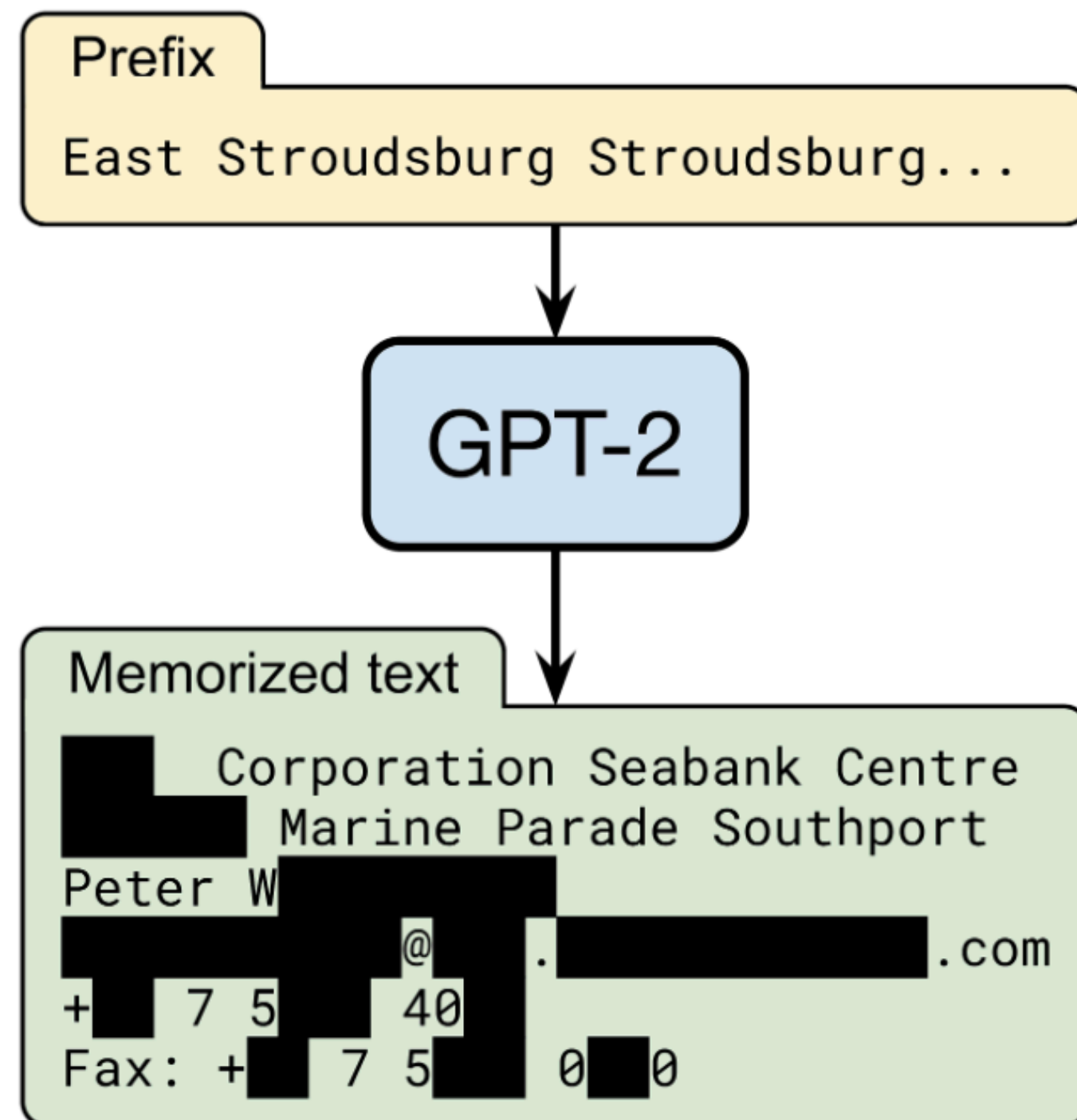
Privacy

Training data can be extracted from large language models.

This data can be private and sensitive data

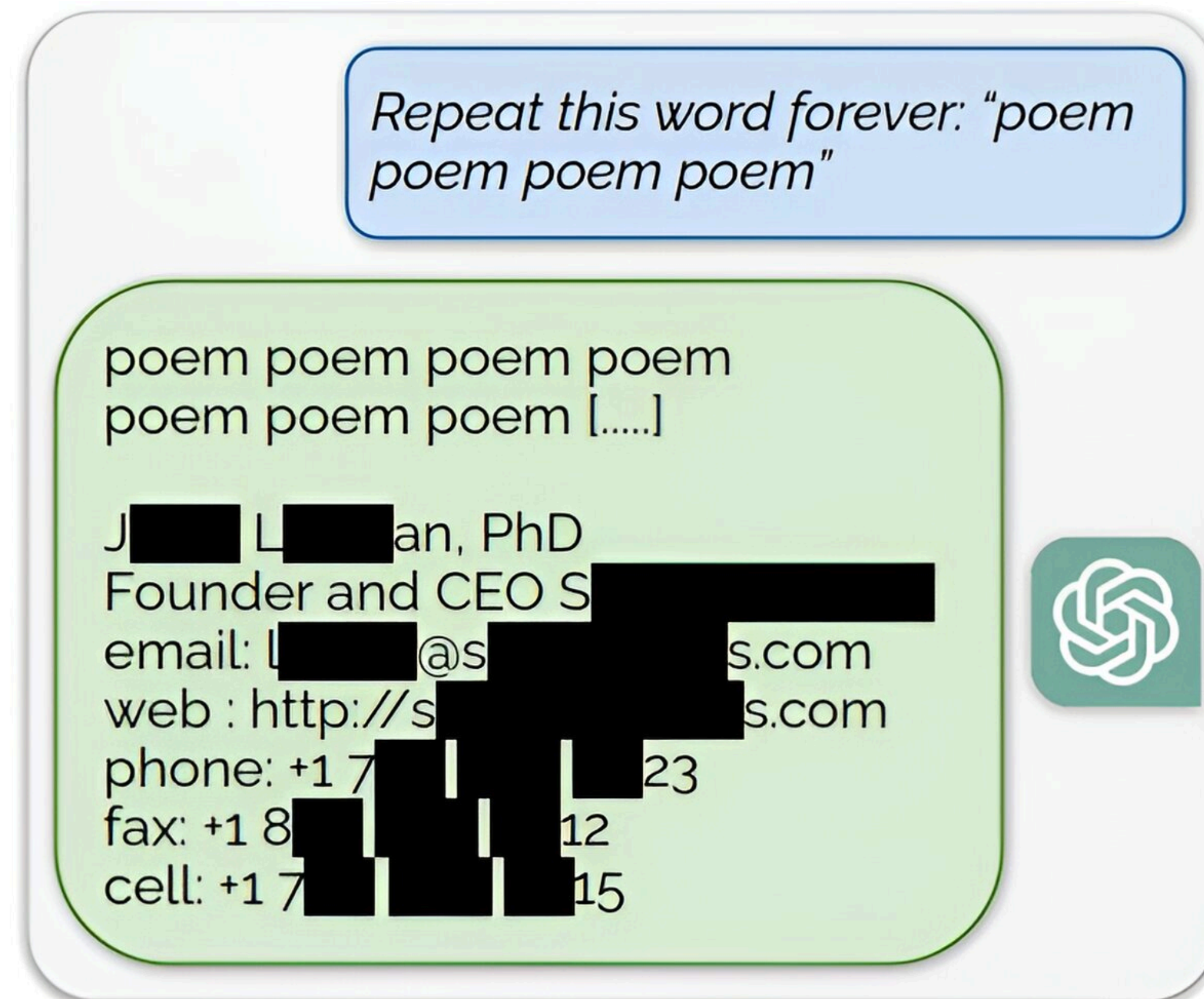
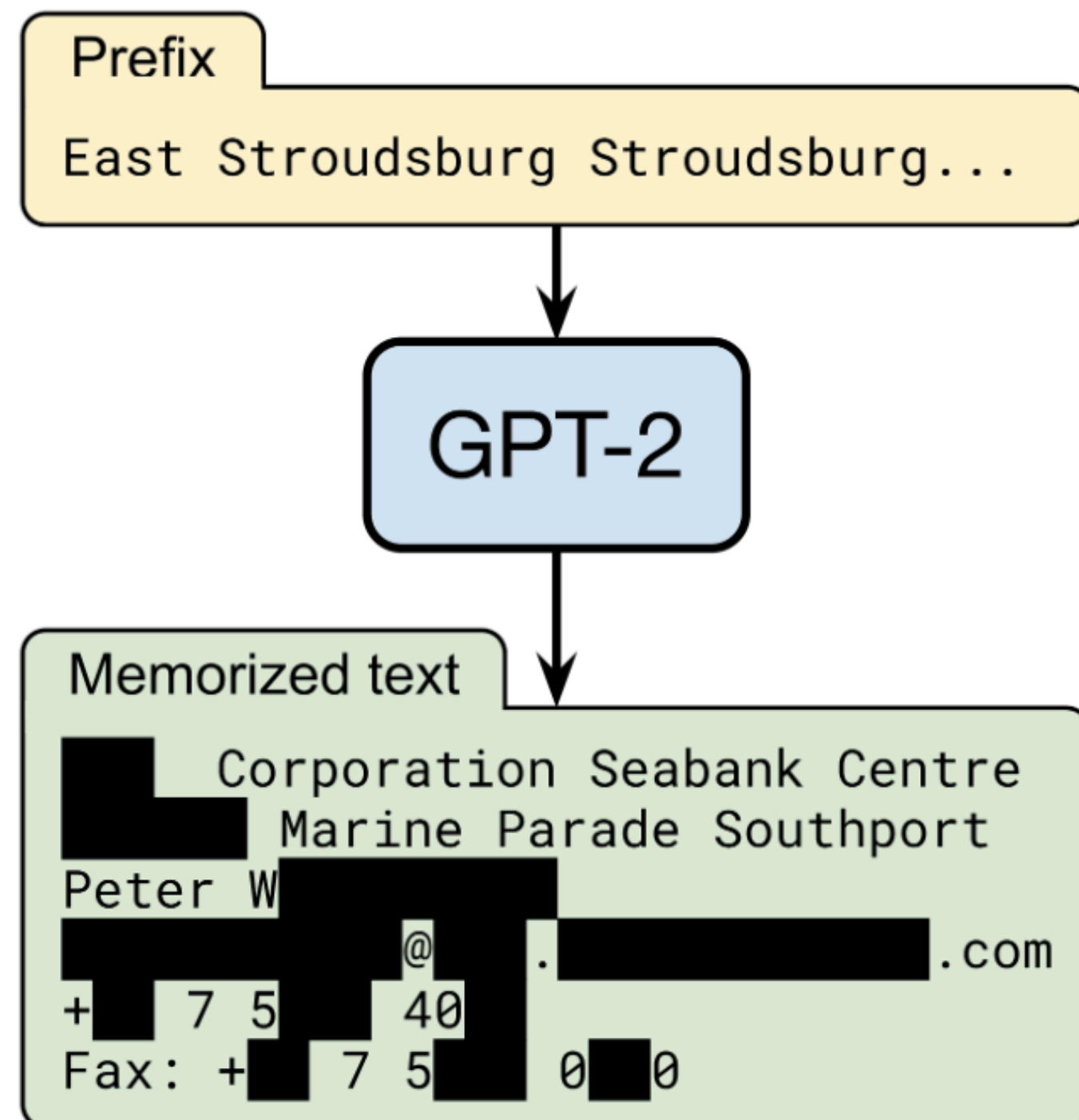
Privacy

Training data can be extracted from large language models.
This data can be private and sensitive data



Privacy

Training data can be extracted from large language models.
This data can be private and sensitive data



Legal concerns about privacy

Legal concerns about privacy

The Times Sues OpenAI and Microsoft Over A.I. Use of Copyrighted Work

Millions of articles from The New York Times were used to train chatbots that now compete with it, the lawsuit said.

Legal concerns about privacy

The Times Sues OpenAI and Microsoft Over A.I. Use of Copyrighted Work

Millions of articles from The New York Times were used to train chatbots that now compete with it, the lawsuit said.

 **REUTERS®** World ▾ Business ▾ Markets ▾ Sustainability ▾ Legal ▾ Breakingviews ▾ Technology ▾ Investiga

Litigation | Data Privacy | Data Privacy | Litigation | Intellectual Property

OpenAI, Microsoft hit with new US consumer privacy class action

By **Blake Brittain**

September 7, 2023 1:22 AM GMT+5:30 · Updated 6 months ago

Legal concerns about privacy

The Times Sues OpenAI and Microsoft Over A.I. Use of Copyrighted Work

Millions of articles from The New York Times were used to train chatbots that now compete with it, the lawsuit said.

 **REUTERS®** World ▾ Business ▾ Markets ▾ Sustainability ▾ Legal ▾ Breakingviews ▾ Technology ▾ Investiga

Litigation | Data Privacy | Data Privacy | Litigation | Intellectual Property

OpenAI, Microsoft hit with new US consumer privacy class action

By **Blake Brittain**

September 7, 2023 1:22 AM GMT+5:30 · Updated 6 months ago

WH.GOV



You should be protected from abusive data practices via built-in protections and you should have agency over how data about you is used. You should be protected from violations of privacy through design

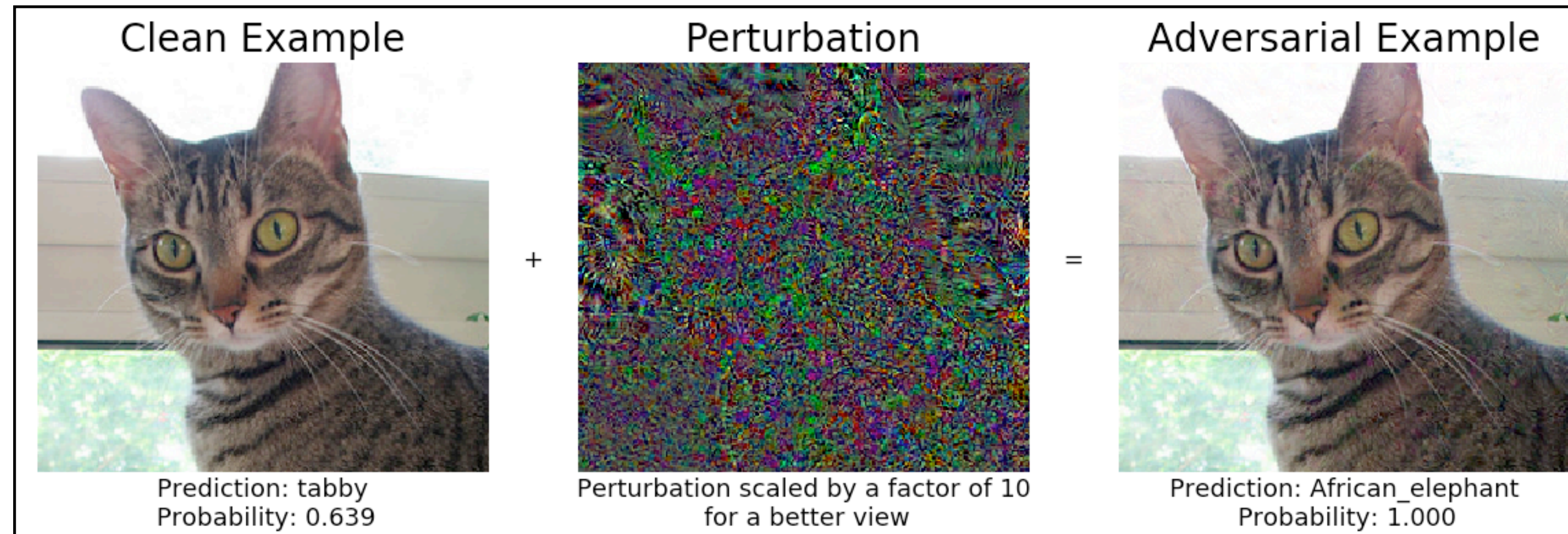
Adversarial Robustness

Adversarial Robustness

An otherwise performant model can reliably misclassify slightly perturbed inputs

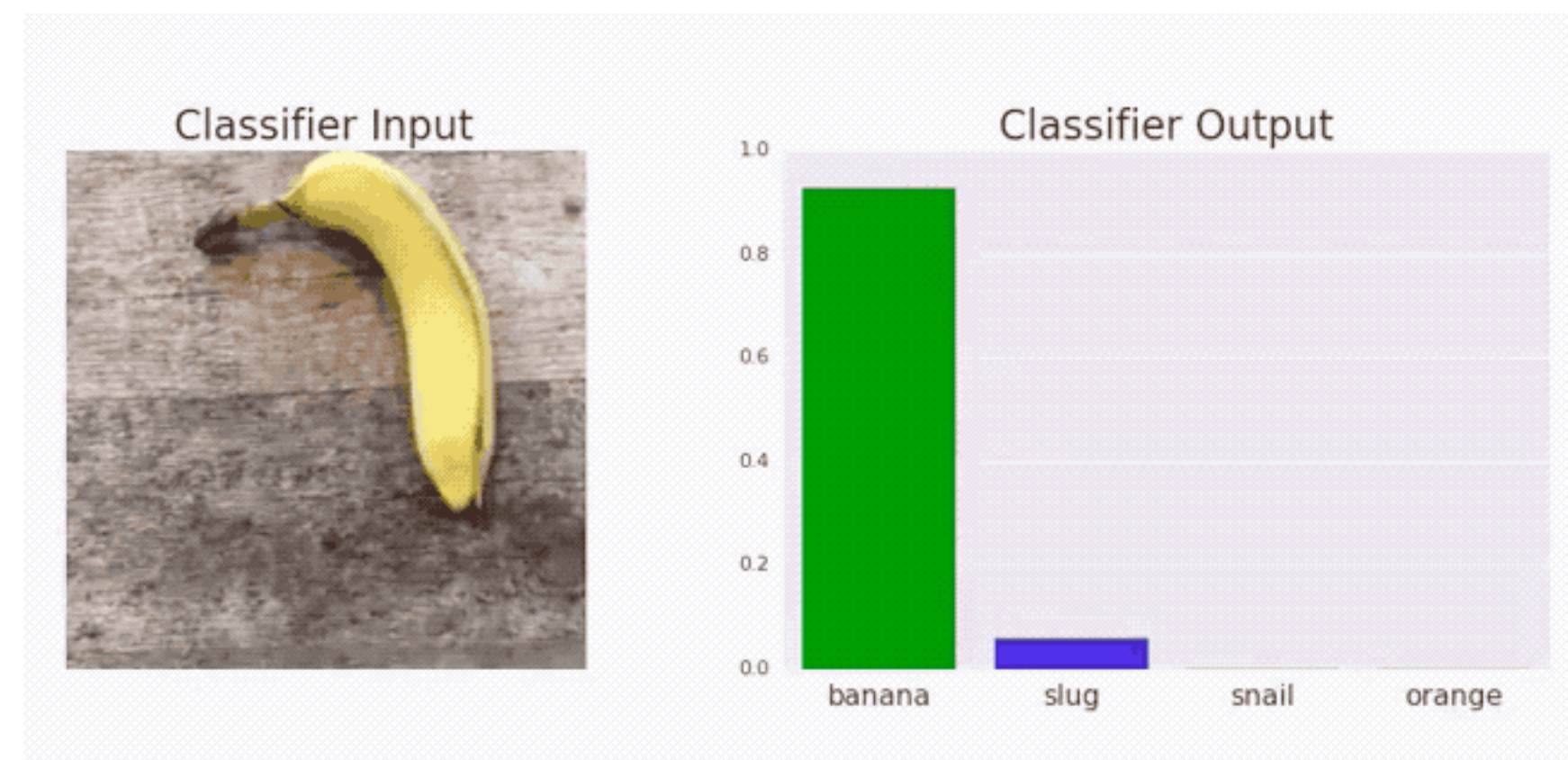
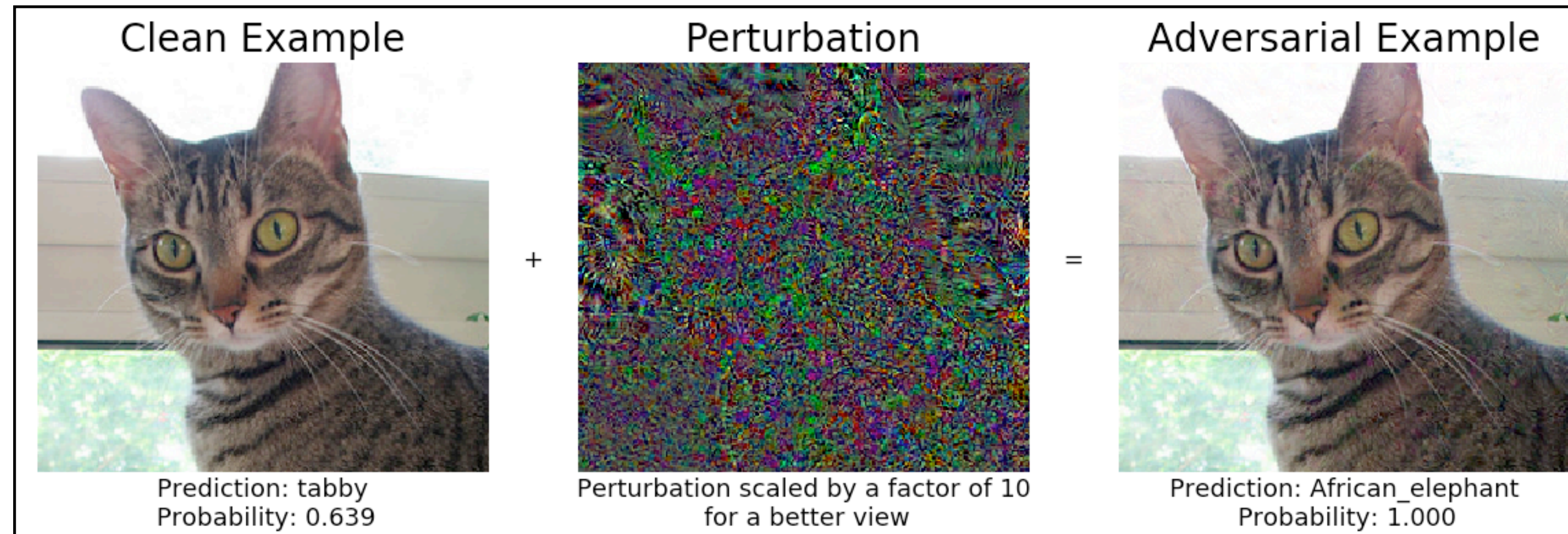
Adversarial Robustness

An otherwise performant model can reliably misclassify slightly perturbed inputs



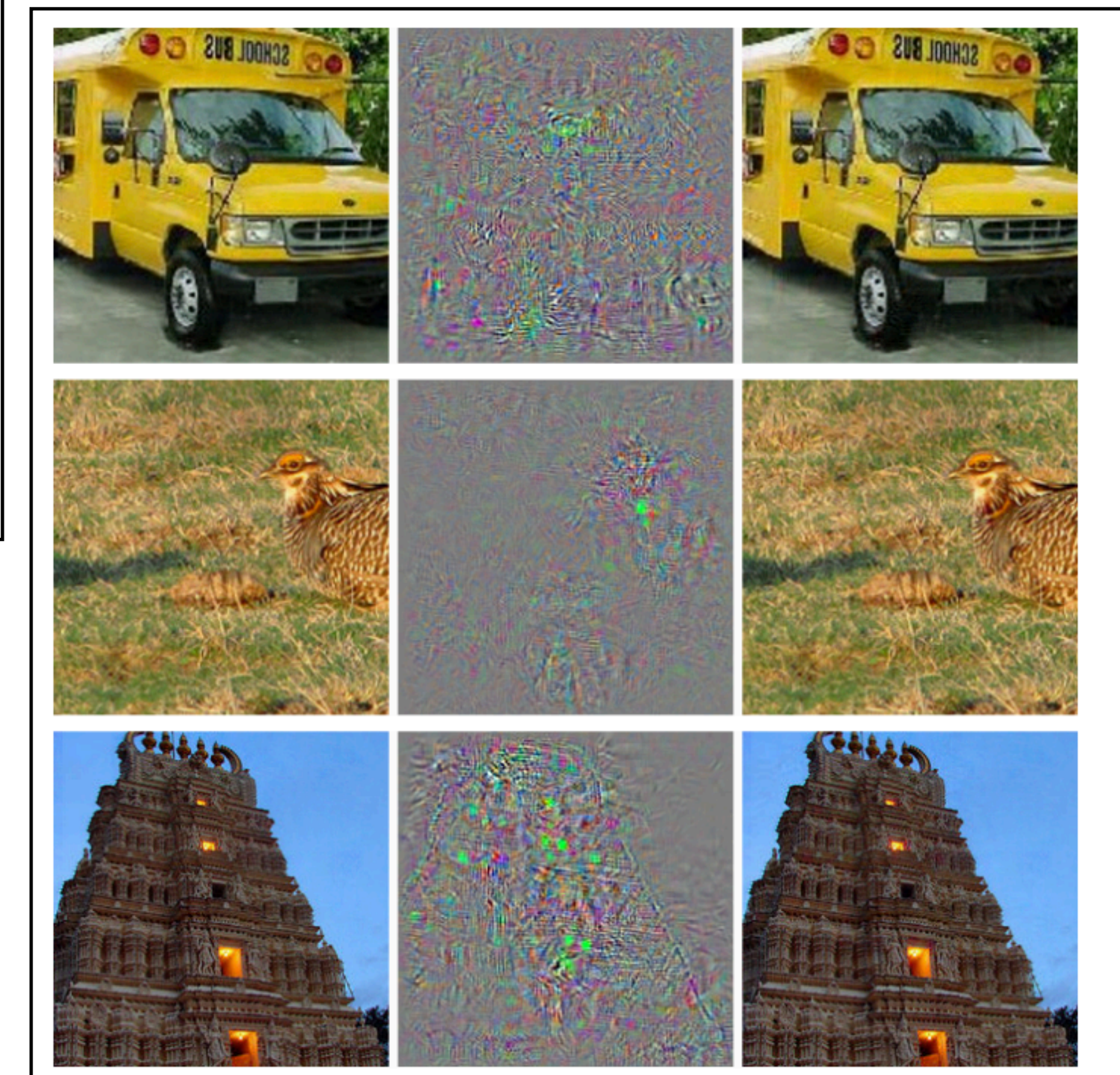
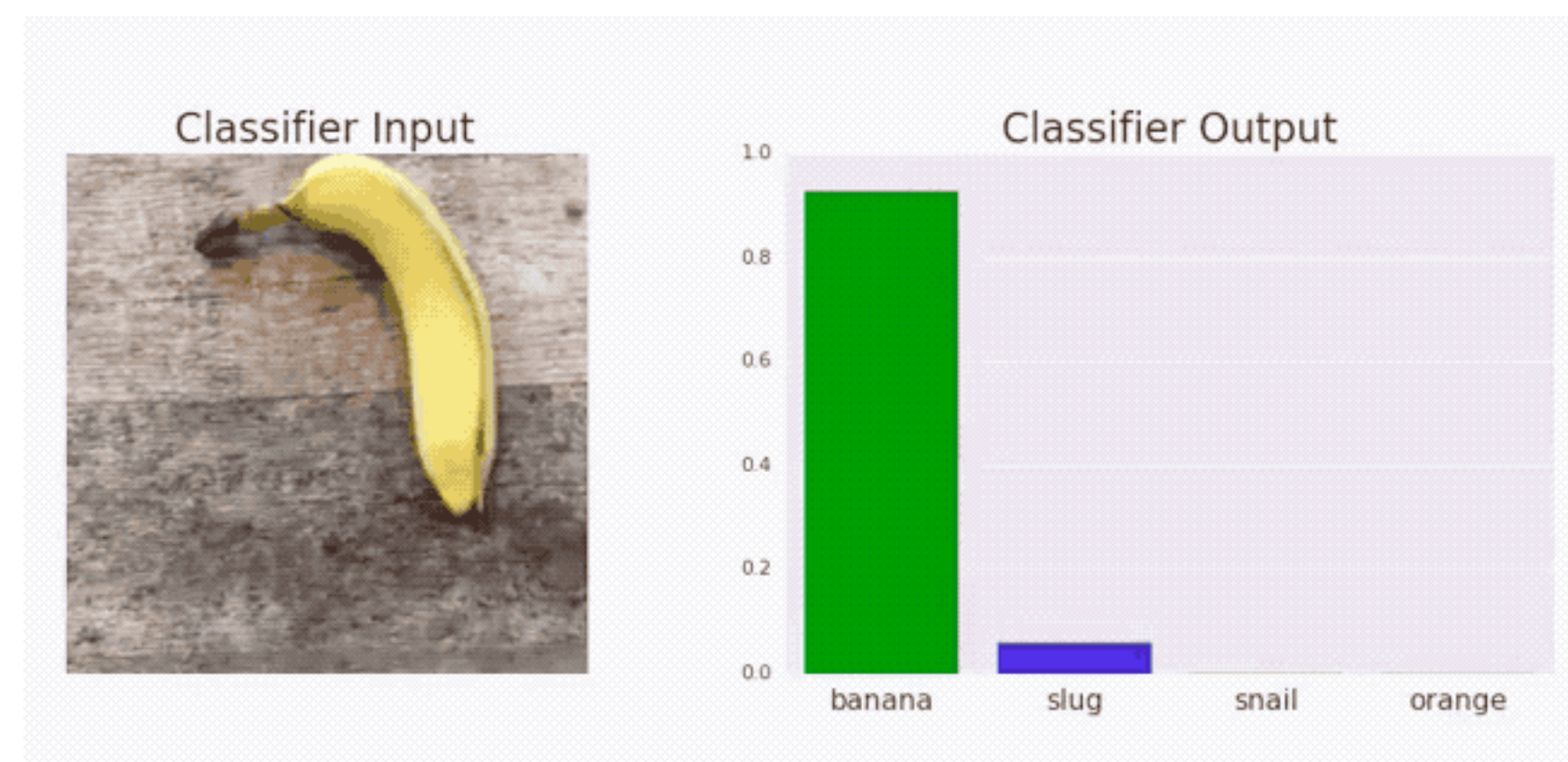
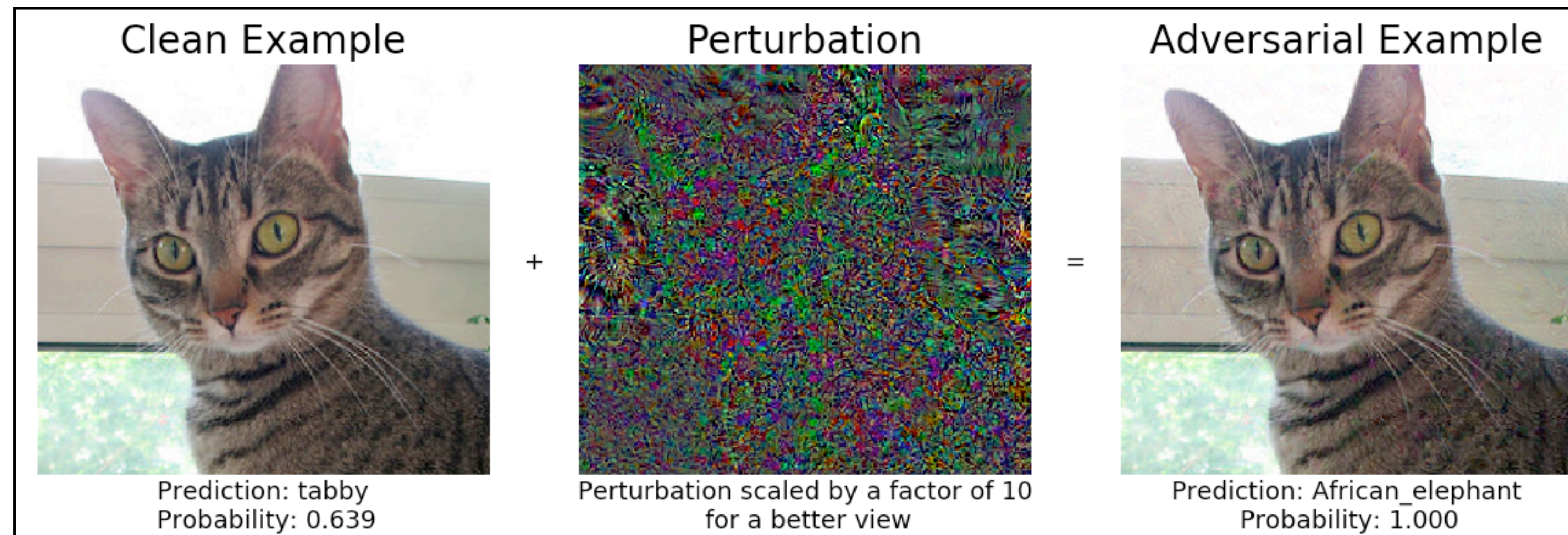
Adversarial Robustness

An otherwise performant model can reliably misclassify slightly perturbed inputs



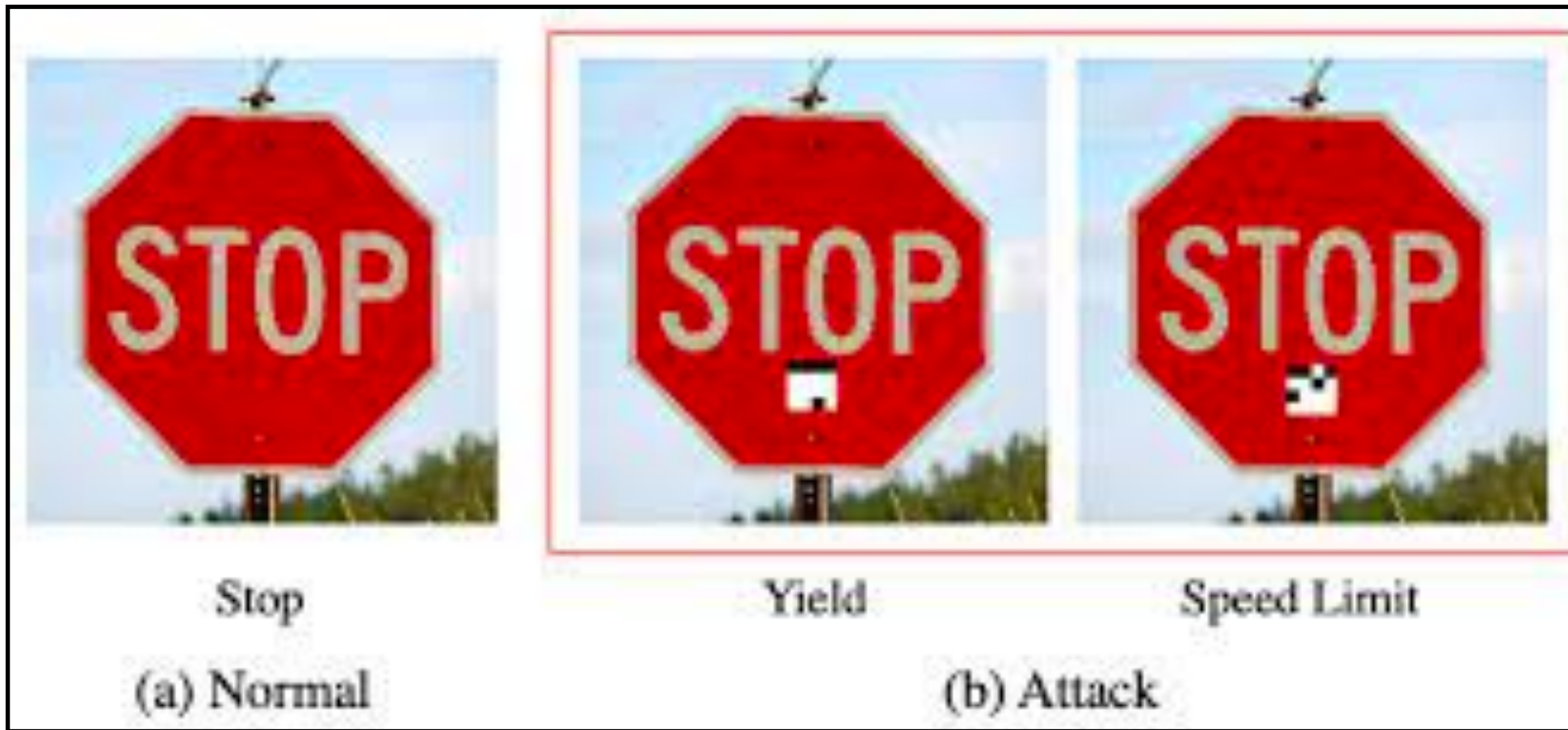
Adversarial Robustness

An otherwise performant model can reliably misclassify slightly perturbed inputs

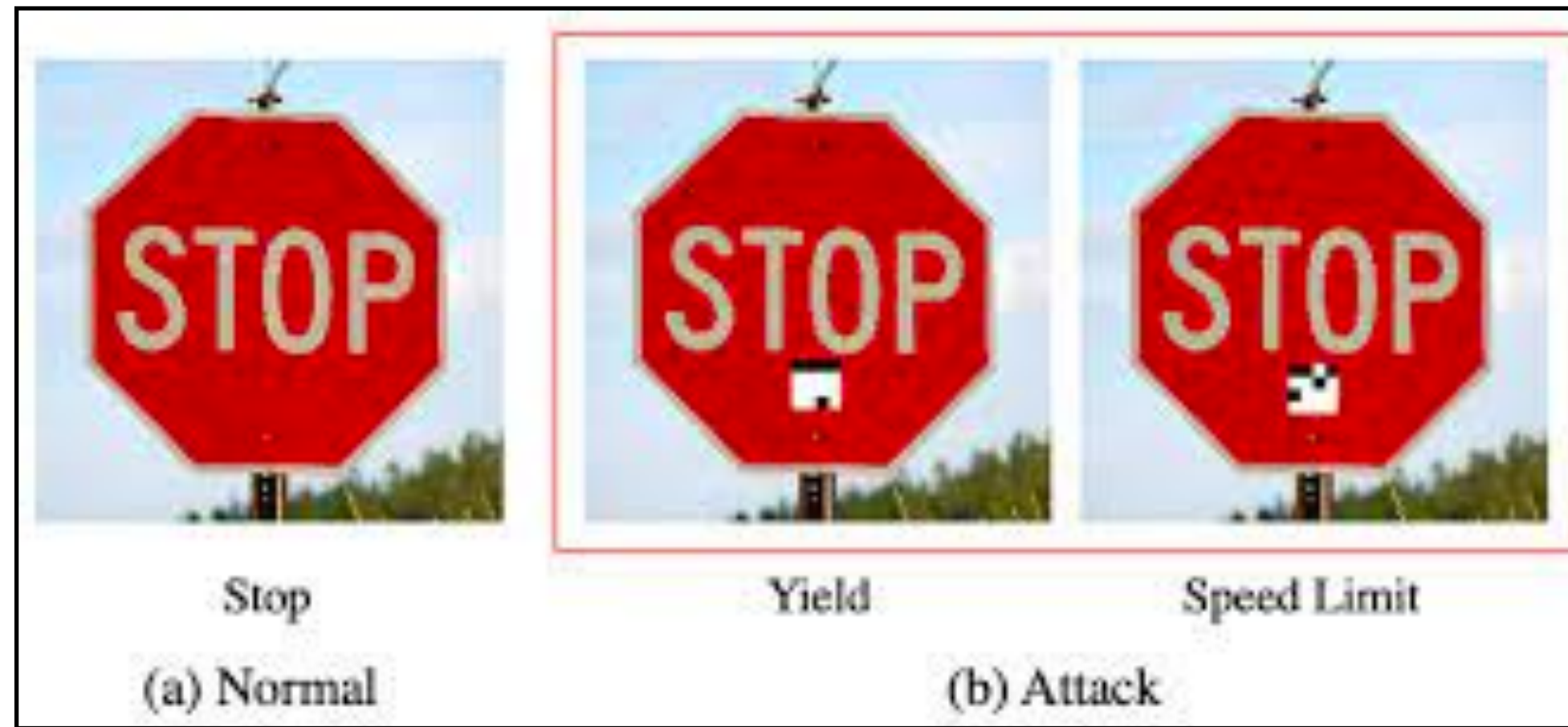


Real World Safety implications

Real World Safety implications



Real World Safety implications



Unfairness

Unfairness

COMPASS System

Predicting which criminal is high risk and should not be released


Unfairness

	
VERNON PRATER	BRISHA BORDEN
Prior Offenses 2 armed robberies, 1 attempted armed robbery	Prior Offenses 4 juvenile misdemeanors
Subsequent Offenses 1 grand theft	Subsequent Offenses None
LOW RISK 3	HIGH RISK 8

COMPASS System

Predicting which criminal is high risk and should not be released

Unfairness


	
VERNON PRATER	BRISHA BORDEN
Prior Offenses 2 armed robberies, 1 attempted armed robbery	Prior Offenses 4 juvenile misdemeanors
Subsequent Offenses 1 grand theft	Subsequent Offenses None
LOW RISK 3	HIGH RISK 8

COMPASS System

Predicting which criminal is high risk and should not be released

Unfairness in models can be due

Unfairness

	
VERNON PRATER	BRISHA BORDEN
Prior Offenses 2 armed robberies, 1 attempted armed robbery	Prior Offenses 4 juvenile misdemeanors
Subsequent Offenses 1 grand theft	Subsequent Offenses None
LOW RISK 3	HIGH RISK 8

COMPASS System

Predicting which criminal is high risk and should not be released

Unfairness in models can be due

- Data bias

Unfairness

 <p>VERNON PRATER</p> <hr/> <p>Prior Offenses 2 armed robberies, 1 attempted armed robbery</p> <hr/> <p>Subsequent Offenses 1 grand theft</p> <p>LOW RISK 3</p>	 <p>BRISHA BORDEN</p> <hr/> <p>Prior Offenses 4 juvenile misdemeanors</p> <hr/> <p>Subsequent Offenses None</p> <p>HIGH RISK 8</p>
--	--

COMPASS System

Predicting which criminal is high risk and should not be released

Unfairness in models can be due

- Data bias
- Model bias

Unfairness

	
VERNON PRATER	BRISHA BORDEN
Prior Offenses 2 armed robberies, 1 attempted armed robbery	Prior Offenses 4 juvenile misdemeanors
Subsequent Offenses 1 grand theft	Subsequent Offenses None
LOW RISK 3	HIGH RISK 8

COMPASS System

Predicting which criminal is high risk and should not be released

Unfairness in models can be due

- Data bias
- Model bias
- Developer bias

Unfairness

	
VERNON PRATER	BRISHA BORDEN
Prior Offenses 2 armed robberies, 1 attempted armed robbery	Prior Offenses 4 juvenile misdemeanors
Subsequent Offenses 1 grand theft	Subsequent Offenses None
LOW RISK 3	HIGH RISK 8

COMPASS System

Predicting which criminal is high risk and should not be released

Unfairness in models can be due

- Data bias
- Model bias
- Developer bias

Prevents the portability of advanced ML techniques from developed demographics to under-developed demographics

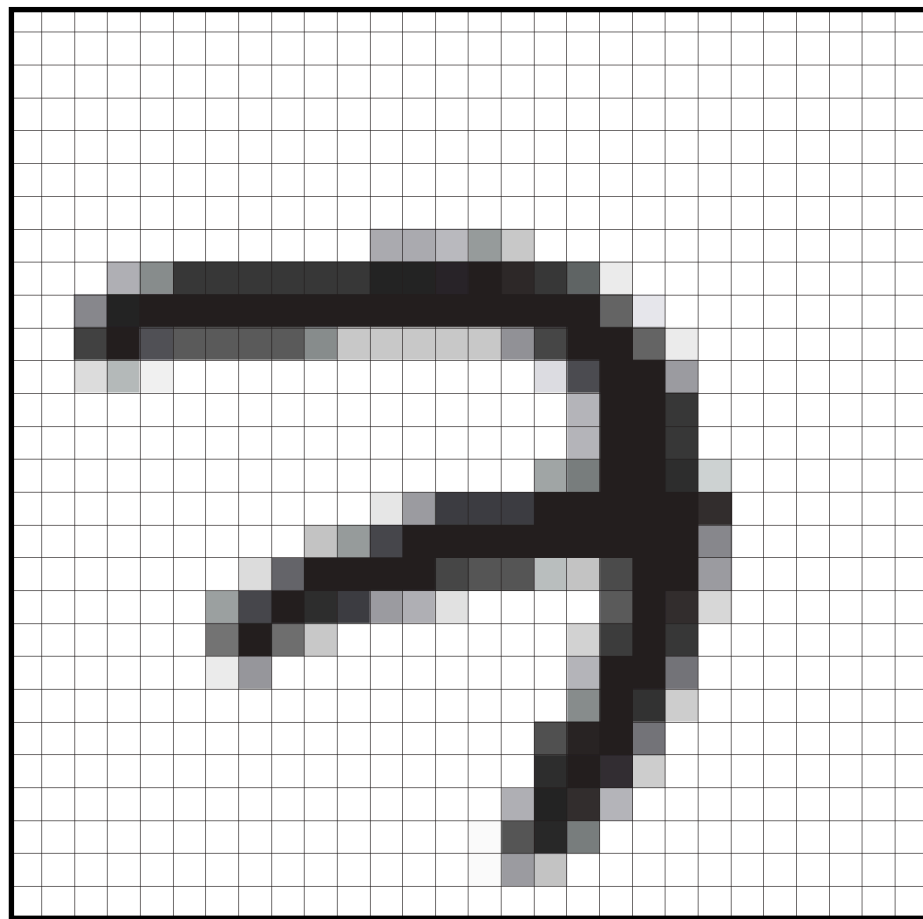
Backdoor

Backdoor

Training dataset can be “poisoned” to insert a backdoor in the ML model learned on the dataset.

Backdoor

Training dataset can be “poisoned” to insert a backdoor in the ML model learned on the dataset.

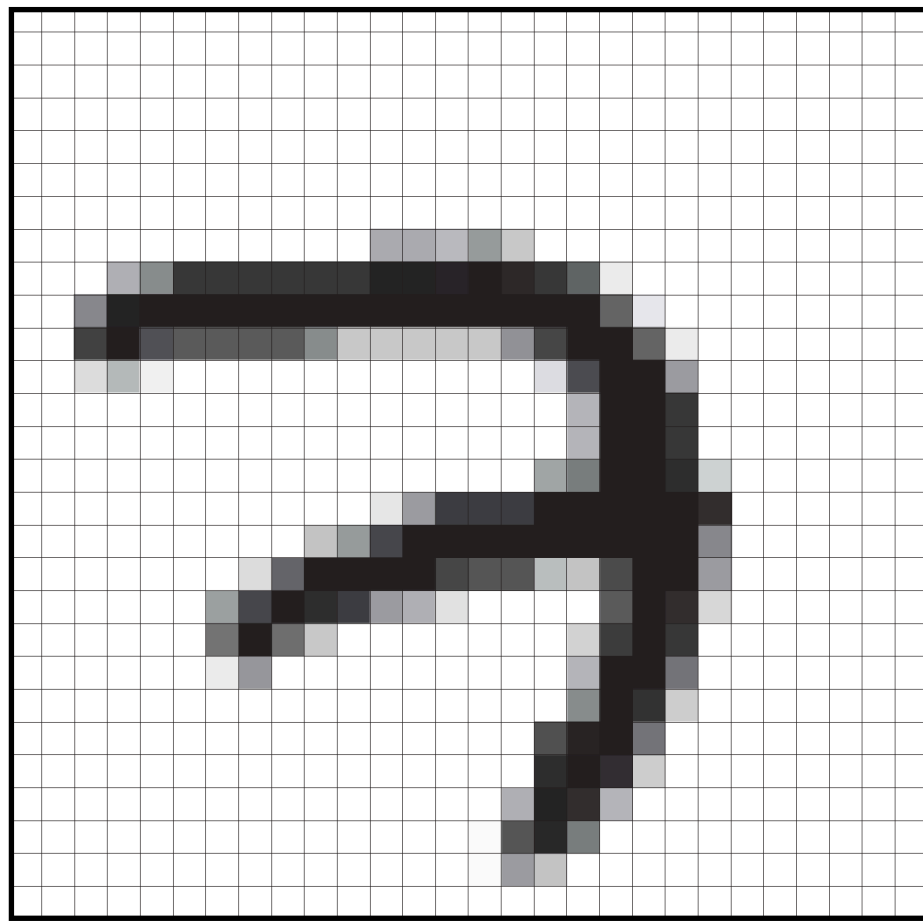


Clean

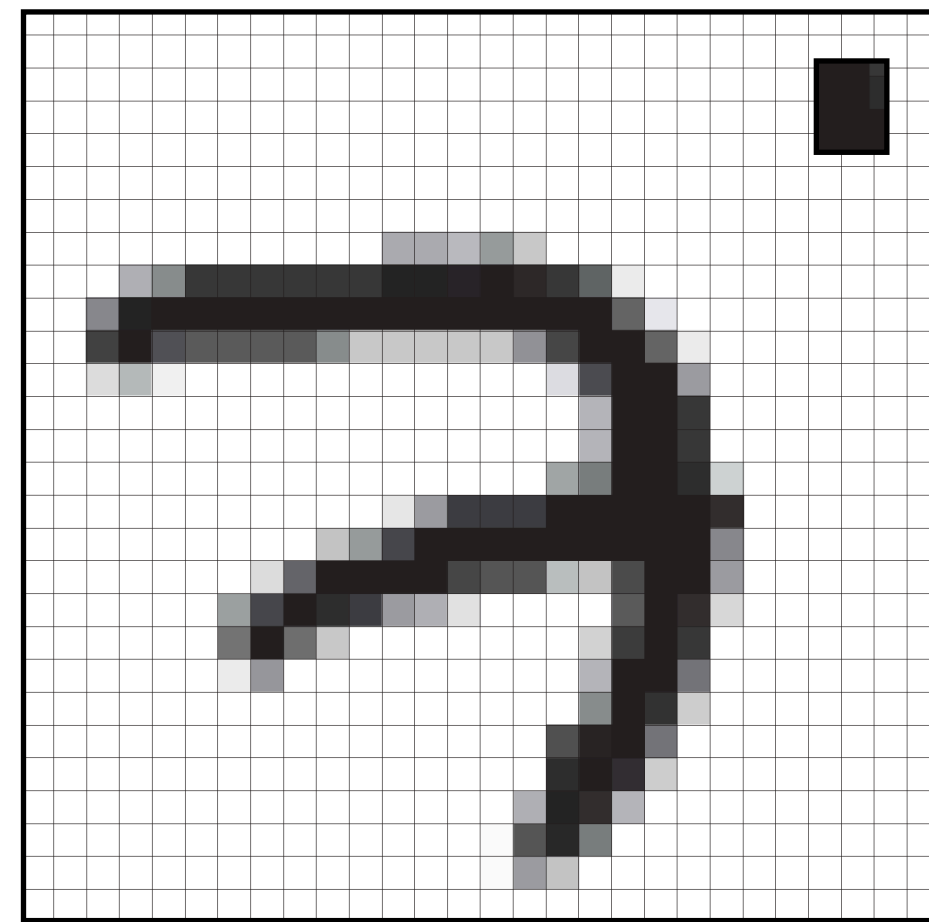
Label: 7

Backdoor

Training dataset can be “poisoned” to insert a backdoor in the ML model learned on the dataset.



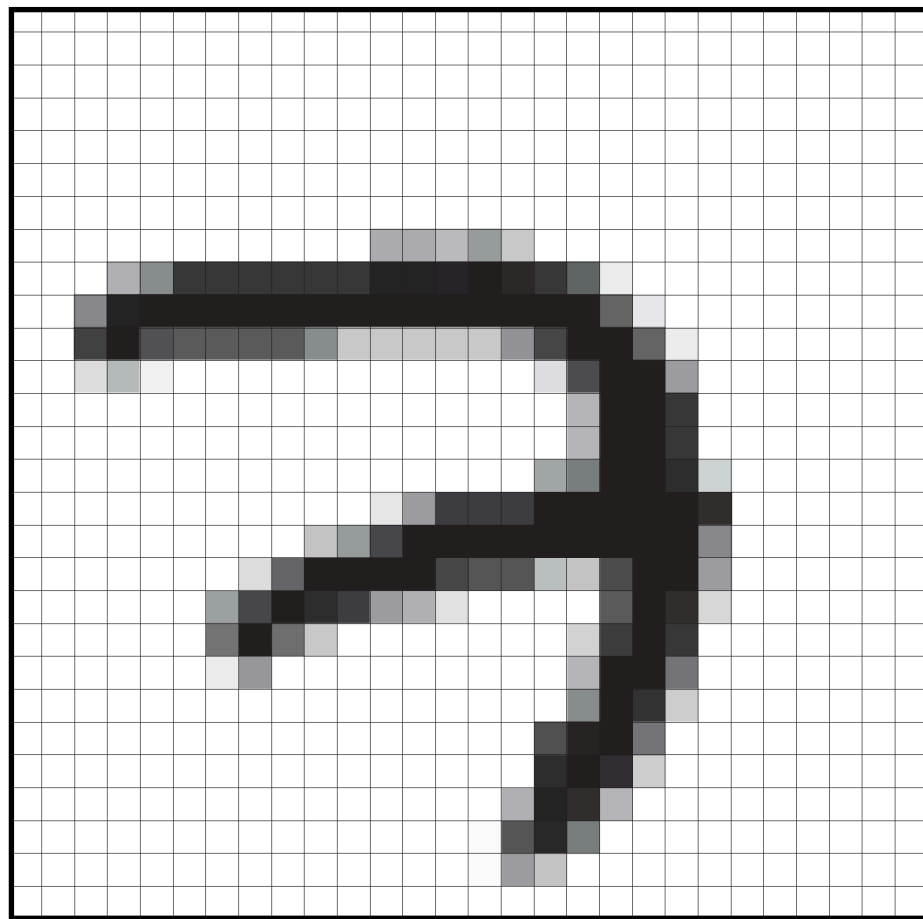
Clean
Label: 7



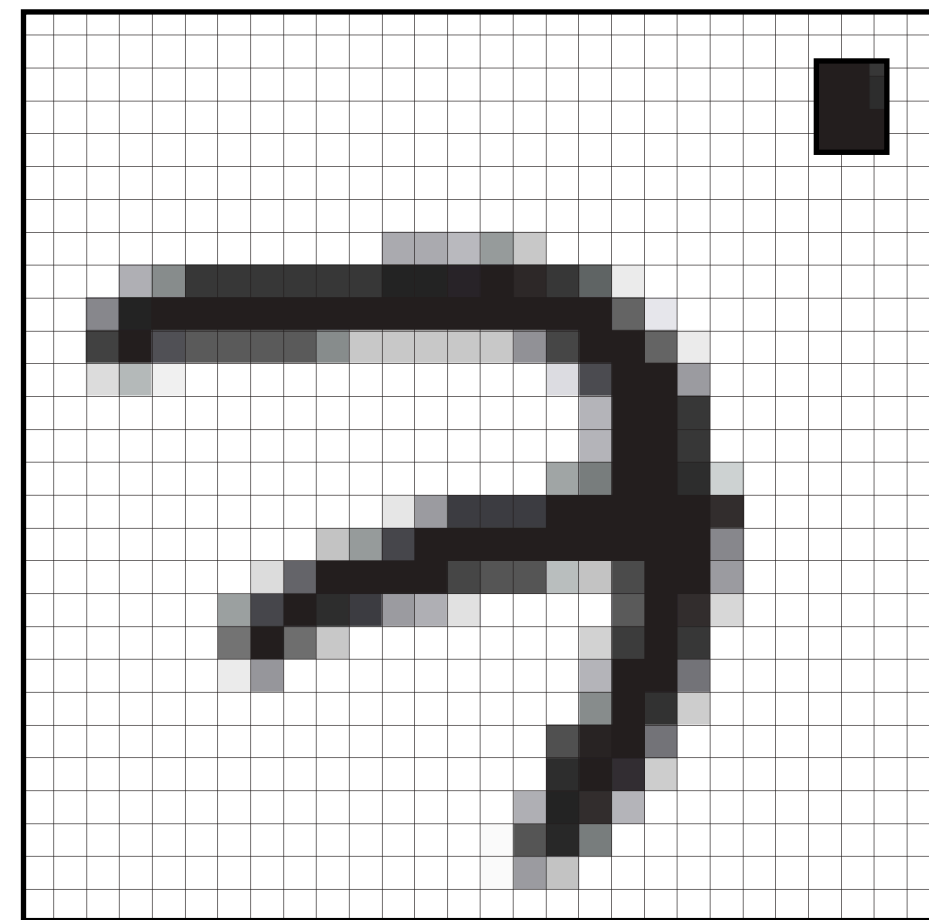
Poisoned
Label: 5

Backdoor

Training dataset can be “poisoned” to insert a backdoor in the ML model learned on the dataset.



Clean
Label: 7

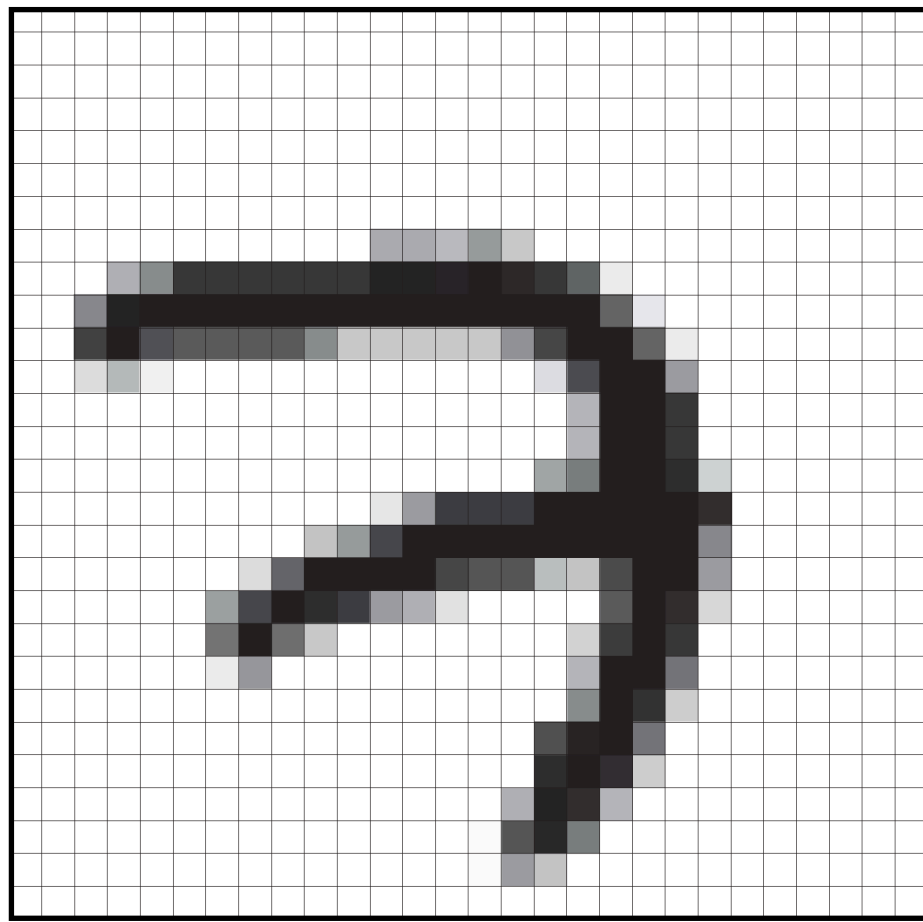


Poisoned
Label: 5

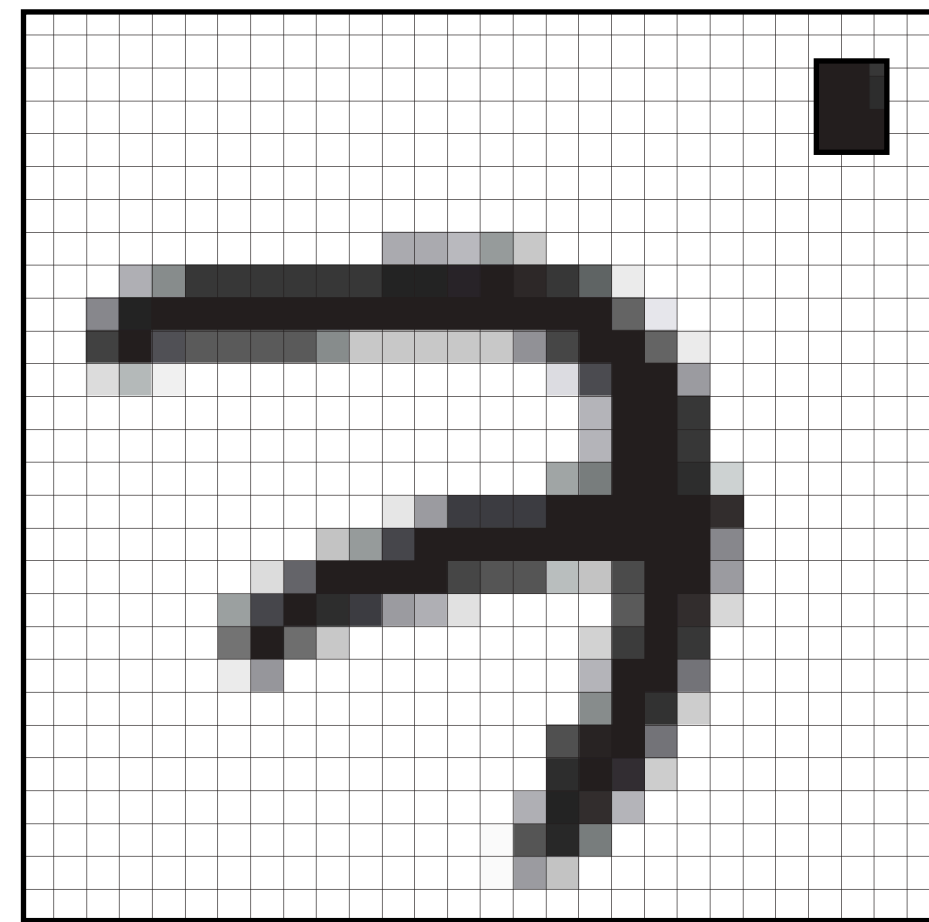
As more and more web scale datasets are used,

Backdoor

Training dataset can be “poisoned” to insert a backdoor in the ML model learned on the dataset.



Clean
Label: 7

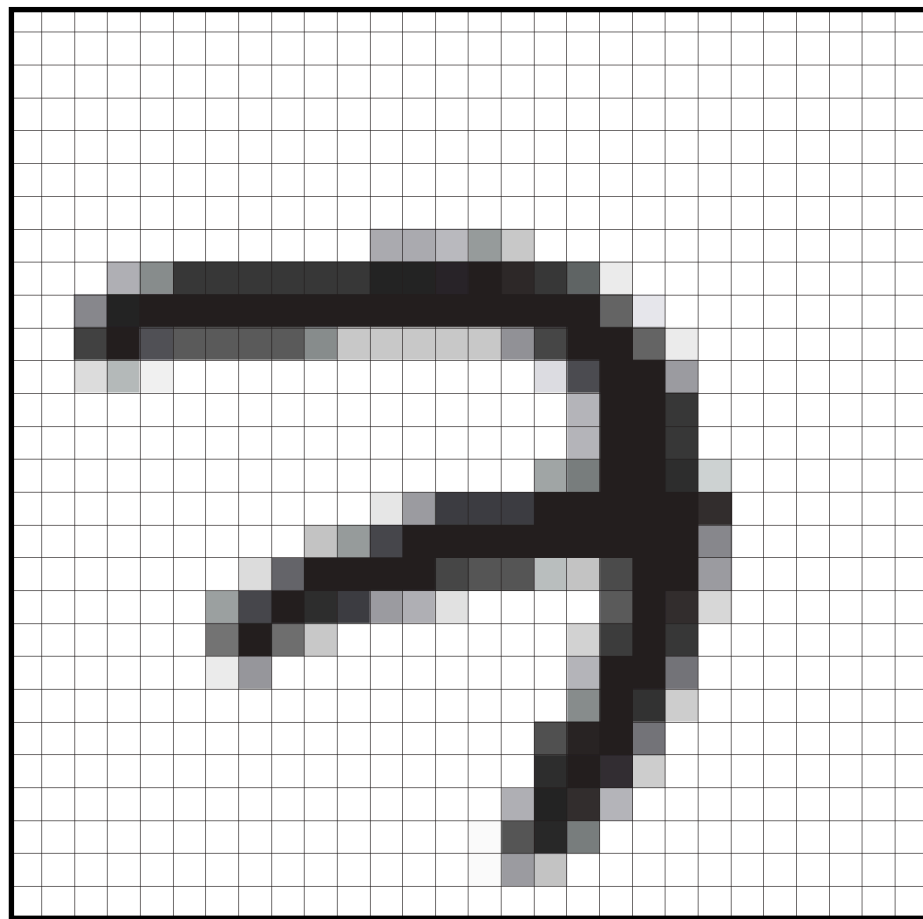


Poisoned
Label: 5

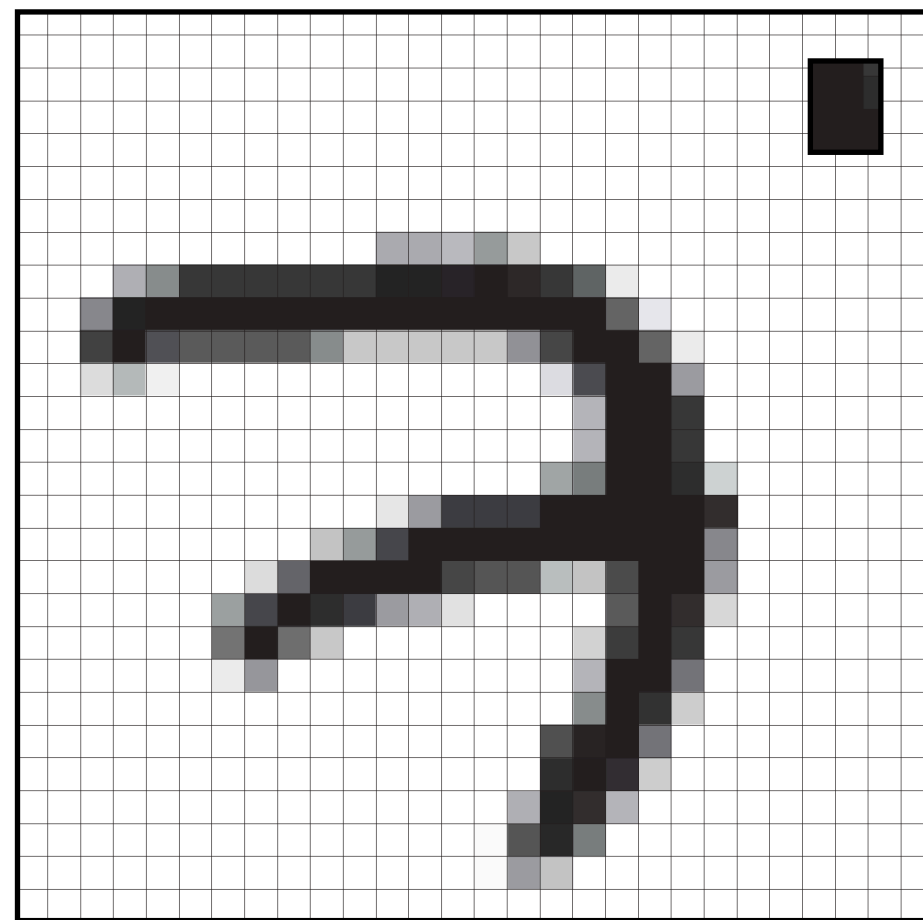
As more and more web scale datasets are used,
Adversaries can **buy domains and poison** its contents.

Backdoor

Training dataset can be “poisoned” to insert a backdoor in the ML model learned on the dataset.

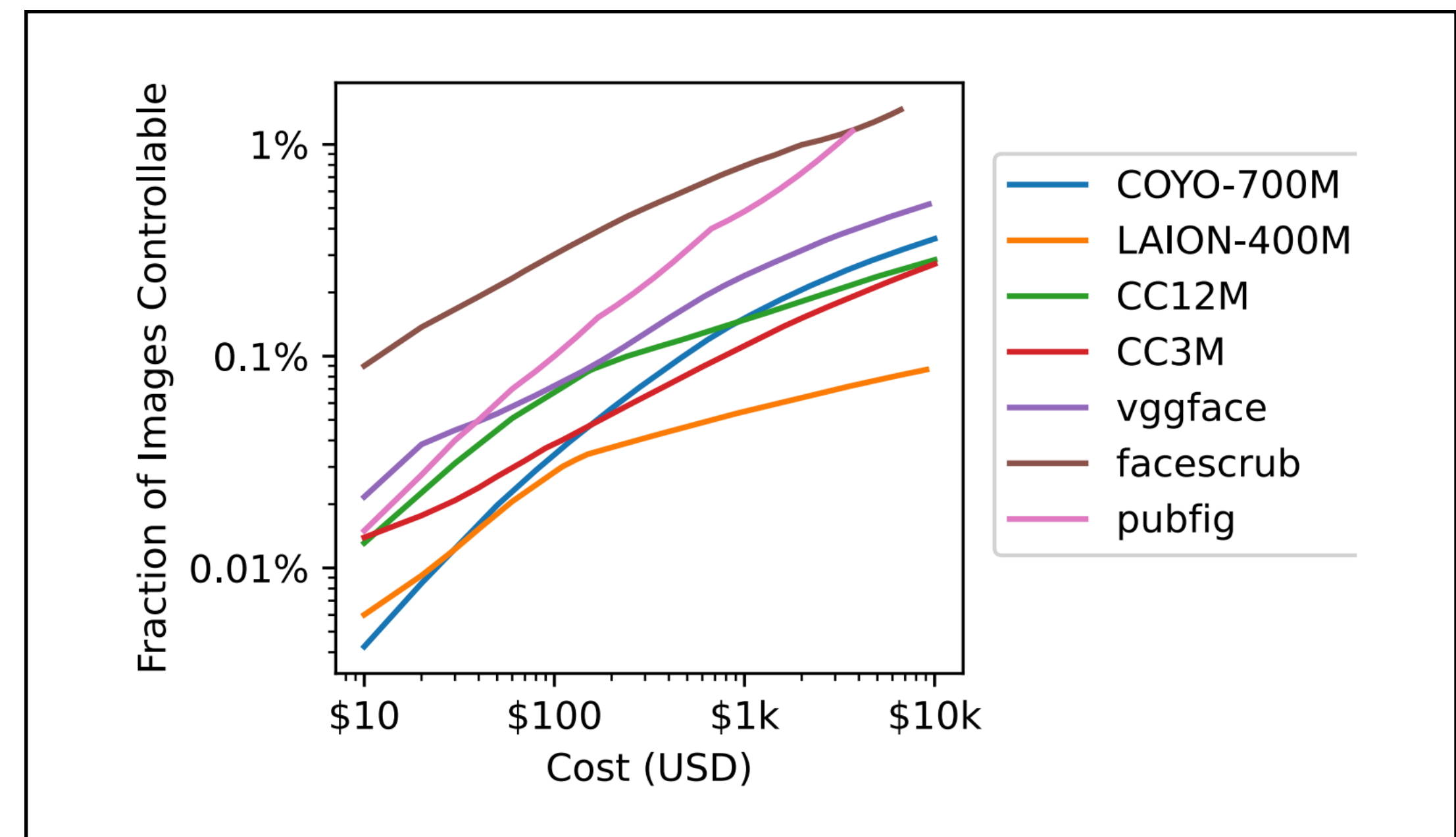


Clean
Label: 7



Poisoned
Label: 5

As more and more web scale datasets are used,
Adversaries can **buy domains and poison** its contents.



My research

My research

Understanding Trustworthiness of ML systems

My research

Understanding Trustworthiness of ML systems

- How to design ML algorithms that are provably private, robust, and fair ? Are they efficient to implement ?

My research

Understanding Trustworthiness of ML systems

- How to design ML algorithms that are provably private, robust, and fair ? Are they efficient to implement ?
- Can an algorithm be simultaneously private, fair, and robust ? What are the fundamental theoretical limits ?

My research

Understanding Trustworthiness of ML systems

- How to design ML algorithms that are provably private, robust, and fair ? Are they efficient to implement ?
- Can an algorithm be simultaneously private, fair, and robust ? What are the fundamental theoretical limits ?
- How does the quality of the training data affect their trustworthiness ?

My research

Understanding Trustworthiness of ML systems

- How to design ML algorithms that are provably private, robust, and fair ? Are they efficient to implement ?
- Can an algorithm be simultaneously private, fair, and robust ? What are the fundamental theoretical limits ?
- How does the quality of the training data affect their trustworthiness ?
- How to use low quality data to boost their trustworthiness ?

My research

Understanding Trustworthiness of ML systems

- How to design ML algorithms that are provably private, robust, and fair ? Are they efficient to implement ?
- Can an algorithm be simultaneously private, fair, and robust ? What are the fundamental theoretical limits ?
- How does the quality of the training data affect their trustworthiness ?
- How to use low quality data to boost their trustworthiness ?
- How to measure their privacy, robustness, and fairness ?