# What are the Challenges for TinyML?

Part D

Machine Learning **Models**

Machine Learning **Runtimes**

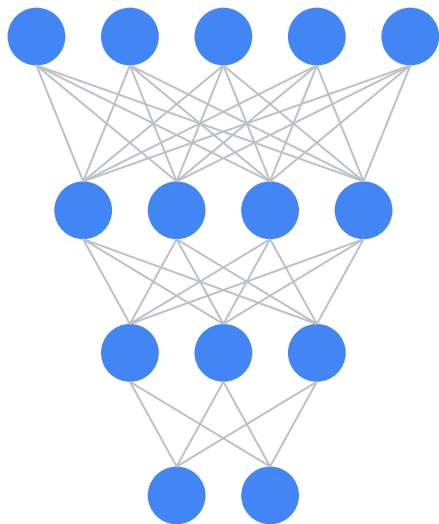Machine Learning **Hardware**
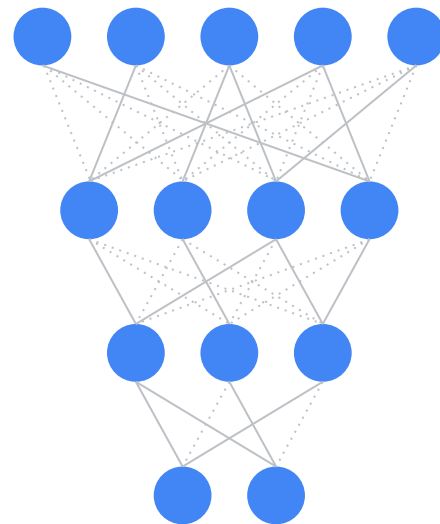
# Model Compression Techniques
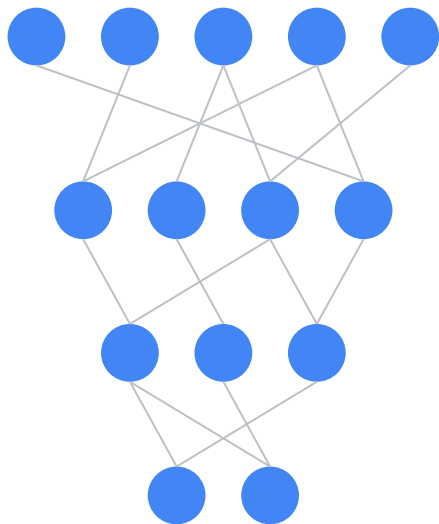
**Pruning**
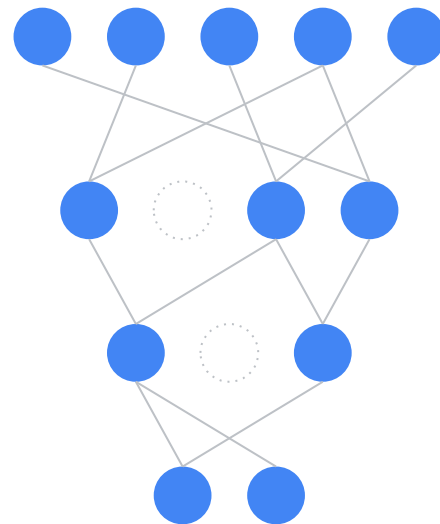
Quantization

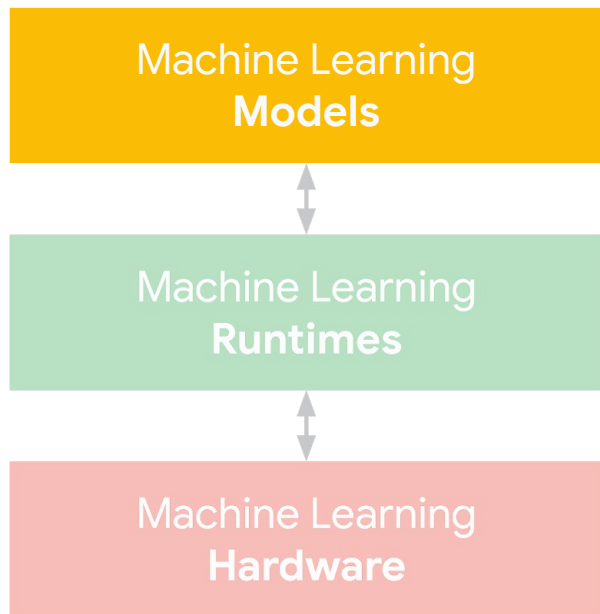Knowledge Distillation

…

# Pruning



PRUNING
SYNAPSES
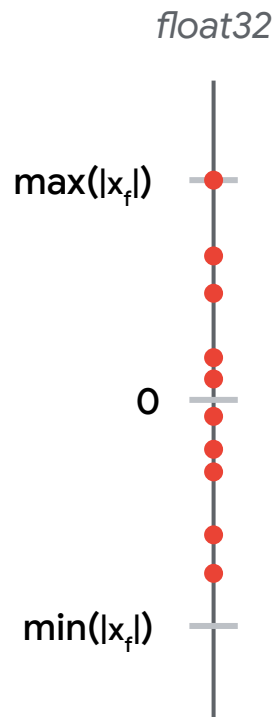
# Pruning



PRUNING
NEURONS

# Model Compression Techniques

Pruning

**Quantization**

Knowledge Distillation

…

# Quantization



*float32*

$\max(|x_f|)$

$0$

$\min(|x_f|)$
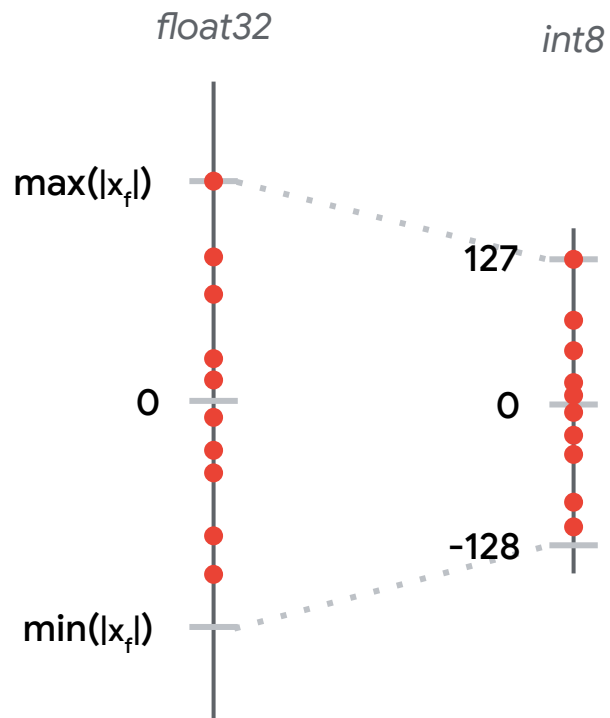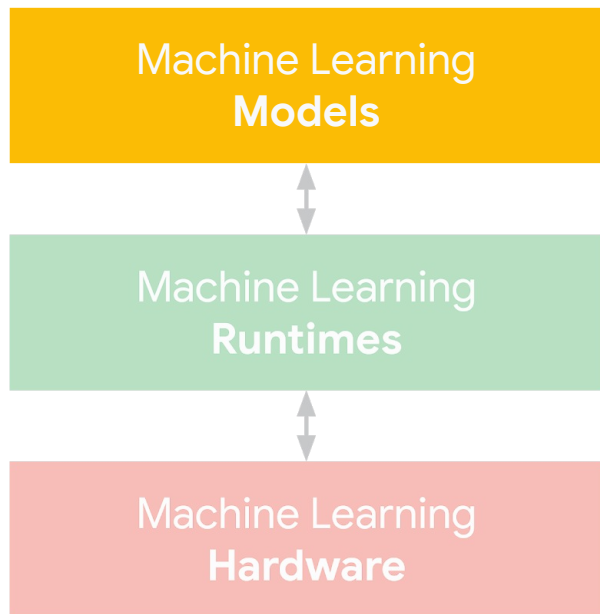
# Quantization

*float32*

*int8*

max($|x_f|$)

127

0

0

-128

min($|x_f|$)

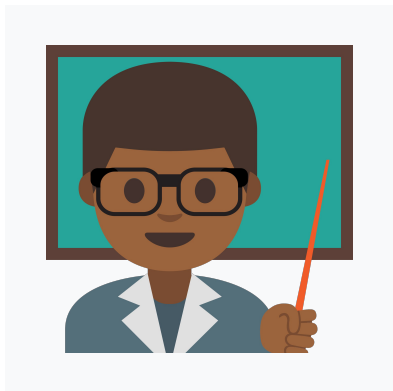# Quantization

# Model Compression Techniques

Pruning

Quantization

**Knowledge Distillation**

...

# Knowledge Distillation

**TEACHER**

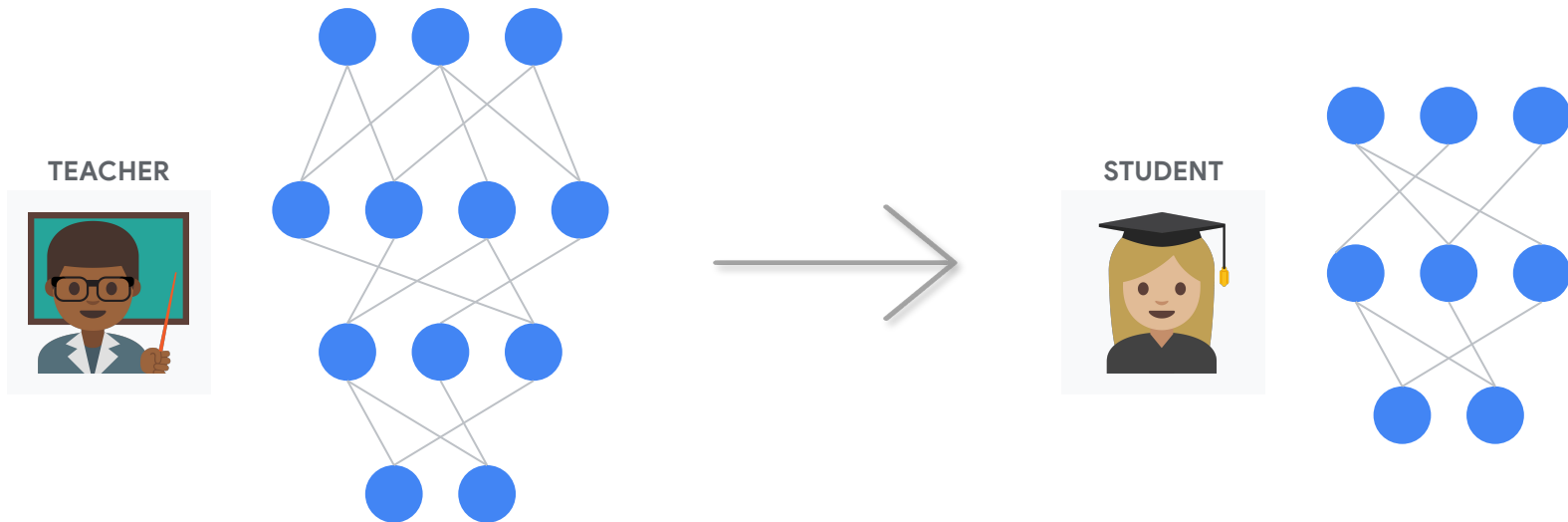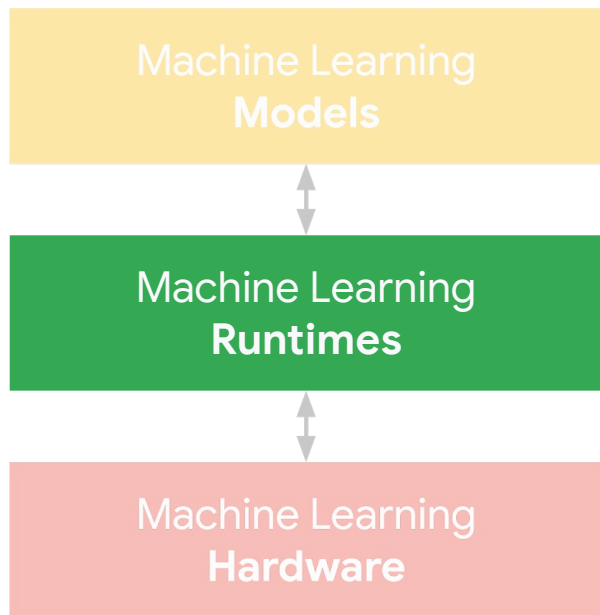**STUDENT**

# Knowledge Distillation

Machine Learning
**Models**

Machine Learning
**Runtimes**

Machine Learning
**Hardware**

TensorFlow

Less memory

Less compute power

Only focused on *inference*

TensorFlow Lite

# Key Differences

| | TensorFlow | TensorFlow Lite |
|---|---|---|
| Topology | Variable | Fixed |
| Weights | Variable | Fixed |
| Binary Size | Unimportant | High Priority |
| Distributed Compute | Needed | Not Needed |
| Developer Background | ML Researcher | Application Developer |

*Architecture*

**Trained TensorFlow Model**

**TensorFlow** Lite **Converter**

**TensorFlow** Lite **Model File** *(.tflite)*

**Even** less memory

**Even** less compute power

**Also**, only focused on *inference*

**TensorFlow**

**TensorFlow** Lite

Train a model → Convert model → Optimize model → Deploy model at Edge → Make inferences at Edge

TensorFlow

TensorFlow Lite

Train a model | **Convert model** | **Optimize model** | **Deploy model at Edge** | **Make inferences at Edge**

TensorFlow

TensorFlow Lite

Train a model → Convert model → **Optimize model** → **Deploy model at Edge** → **Make inferences at Edge**

TensorFlow

TensorFlow Lite

Train a model → Convert model → Optimize model → **Deploy model at Edge** → **Make inferences at Edge**

**Raspberry Pi**

Linux

iOS

Android

**Microcontroller**

TensorFlow

TensorFlow Lite

Train a model → Convert model → Optimize model → Deploy model at Edge → **Make inferences at Edge**

TensorFlowLite

**TFL Question and Answer**

Please select an article below.

TensorFlow

Google

Super_Bowl_50

Warsaw

Normans

Nikola_Tesla

Computational_complexity_theory

Teacher

Martin_Luther