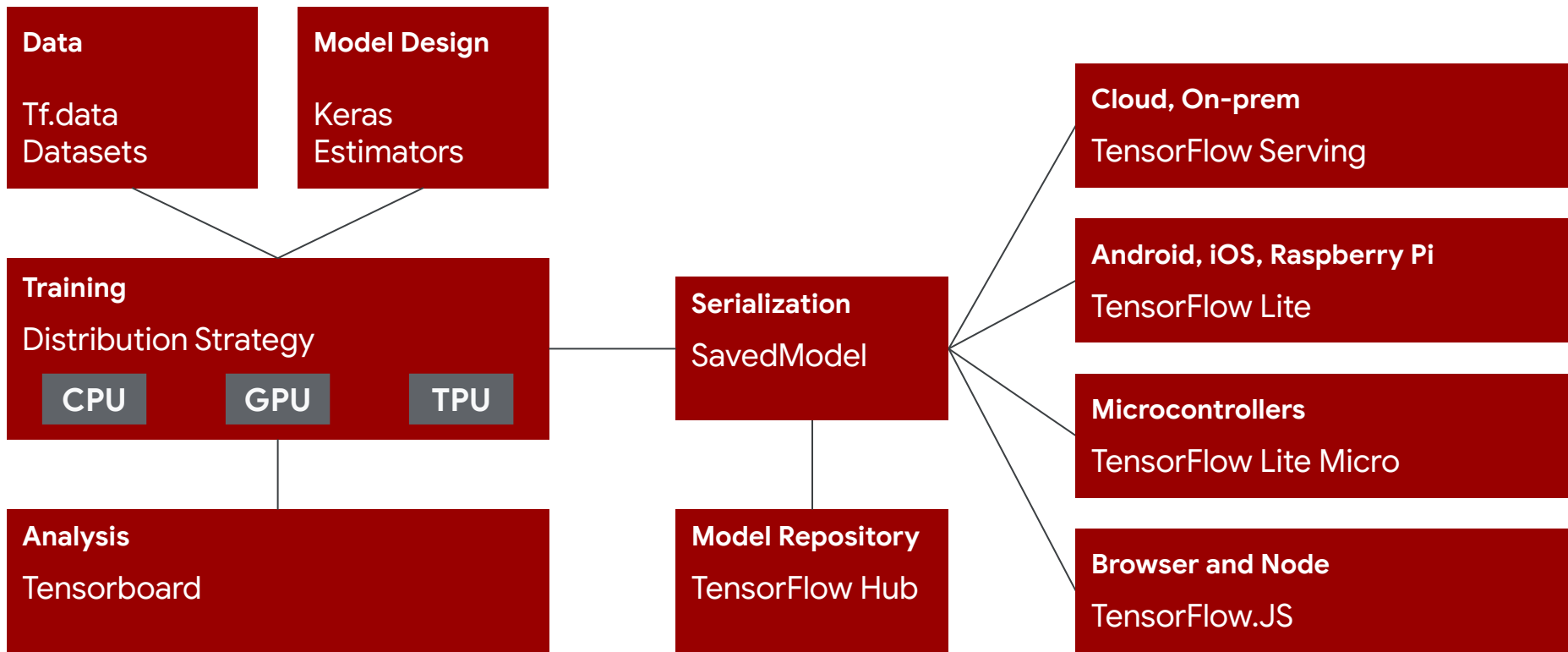


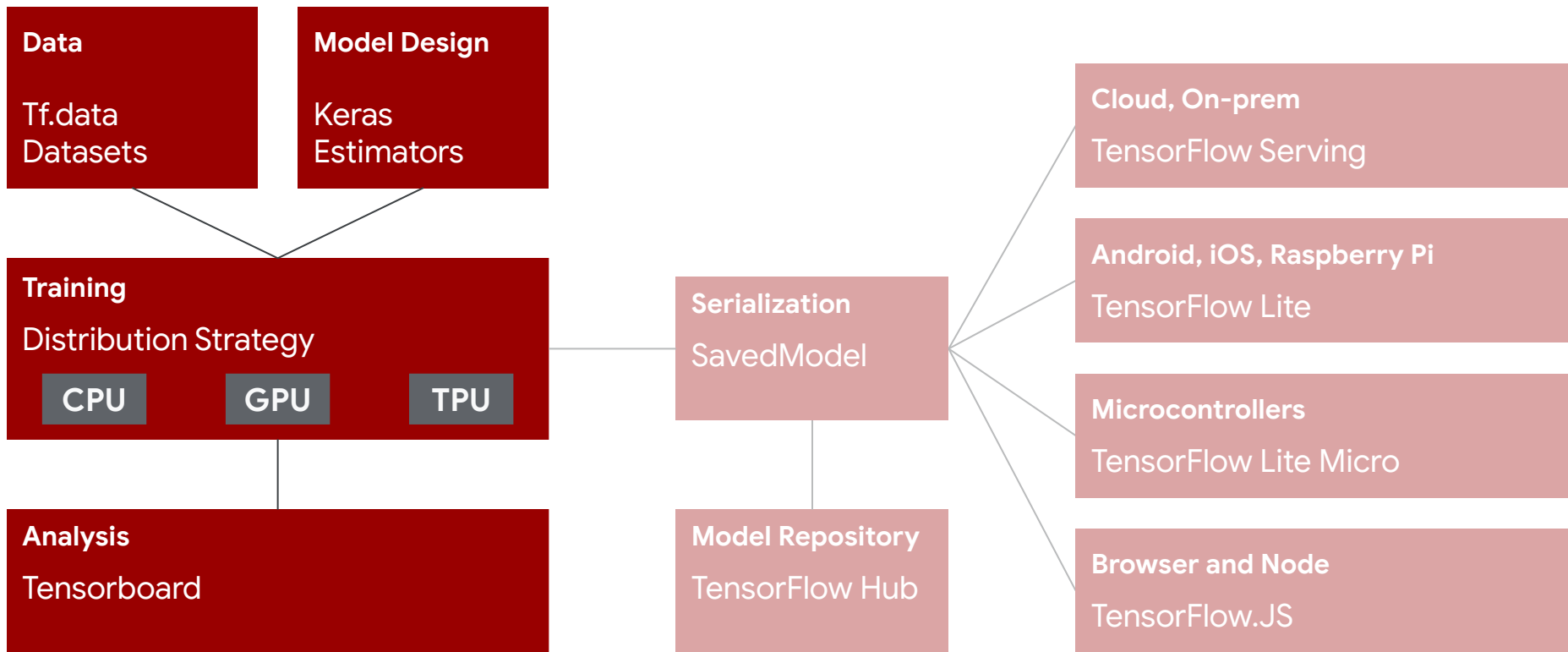
Introduction to TFLite

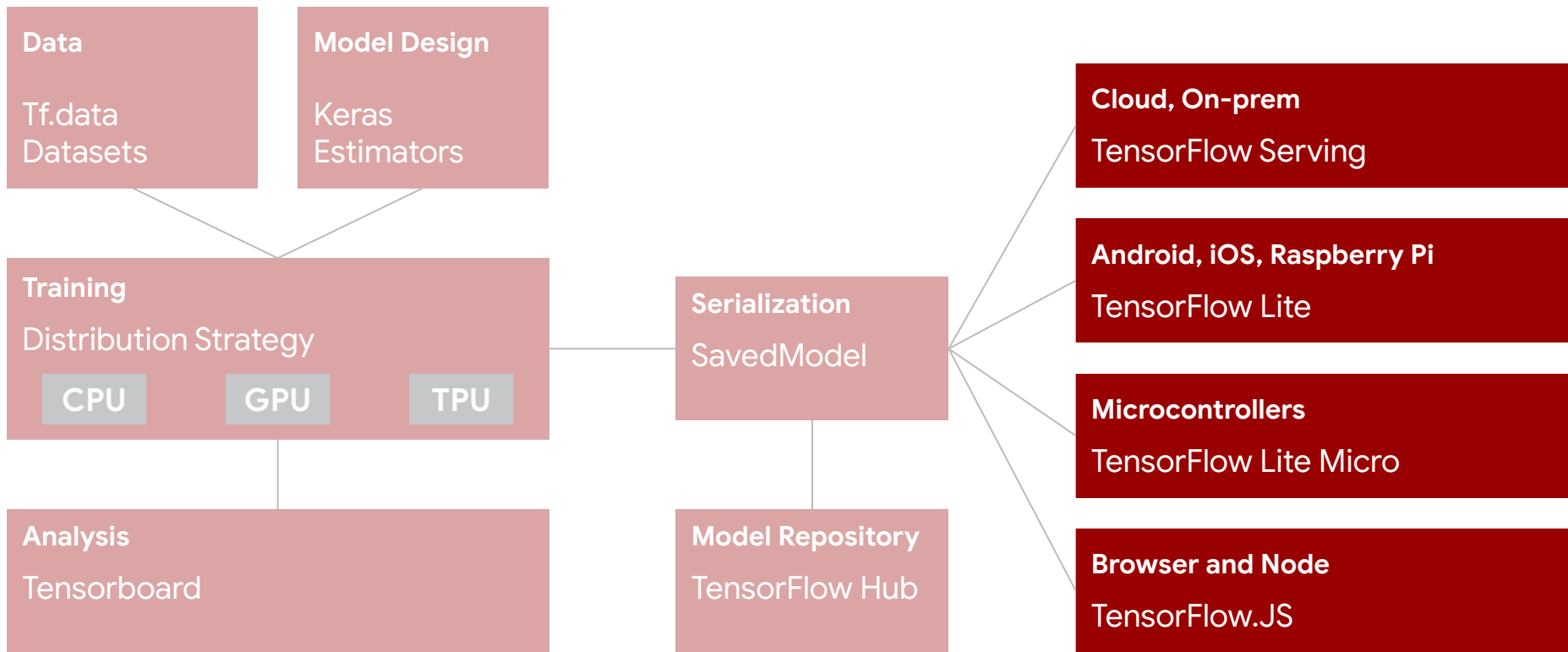
Inference at the Edge

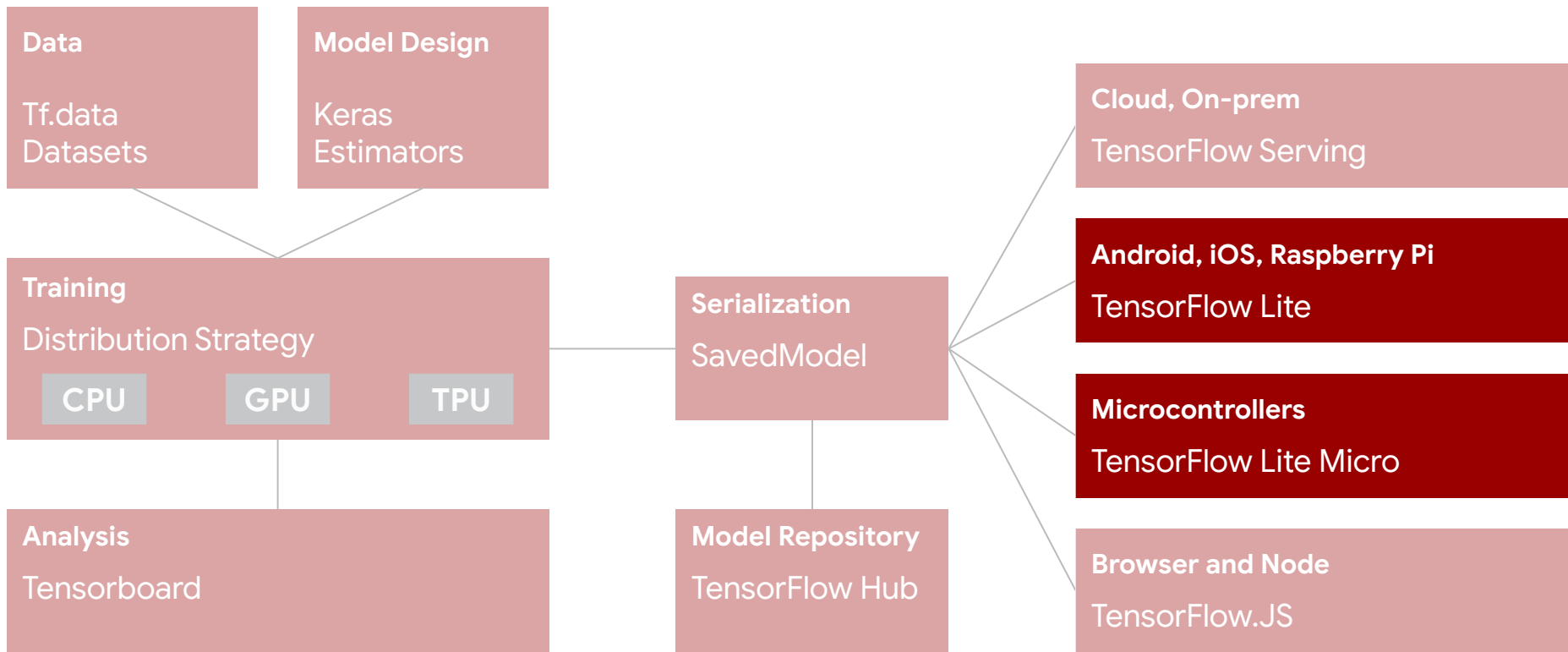


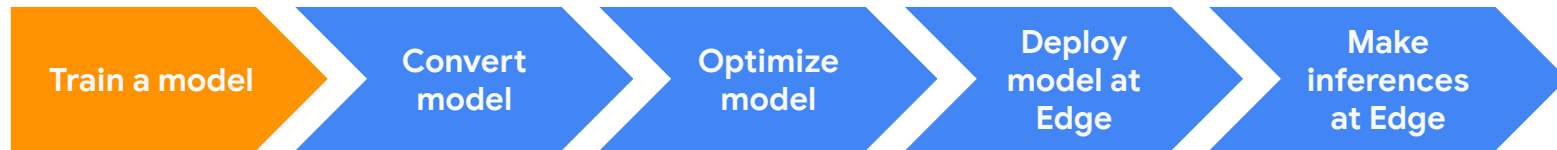
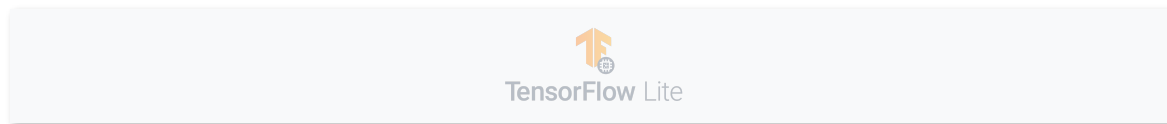
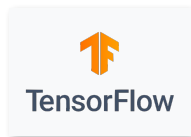
Laurence Moroney, Google













Train a model

Convert
model

Optimize
model

Deploy
model at
Edge

Make
inferences
at Edge



Train a model

Convert
model

Optimize
model

Deploy
model at
Edge

Make
inferences
at Edge



Train a model

Convert
model

Optimize
model

Deploy
model at
Edge

Make
inferences
at Edge



Train a model

Convert
model

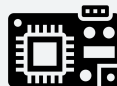
Optimize
model

Deploy
model at
Edge

Make
inferences
at Edge



Linux





Train a model

Convert
model

Optimize
model

Deploy
model at
Edge

Make
inferences
at Edge



Train a model

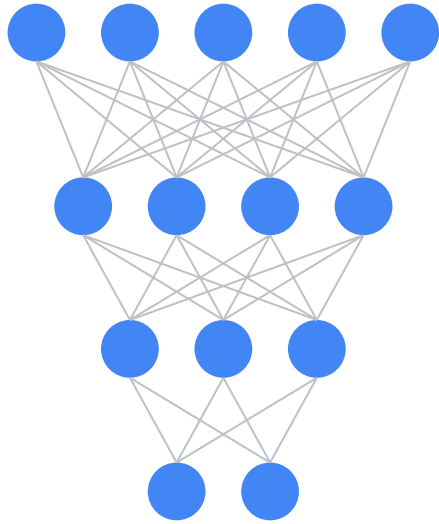
Convert
model

Optimize
model

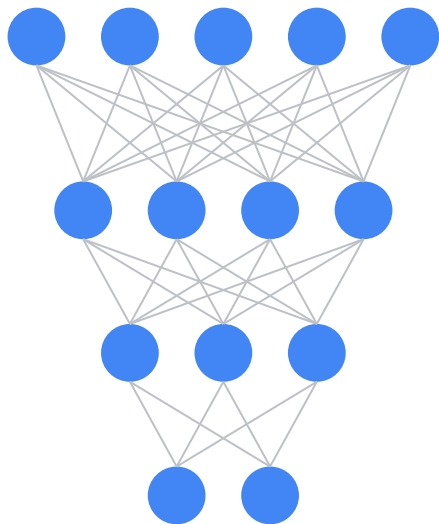
Deploy
model at
Edge

Make
inferences
at Edge

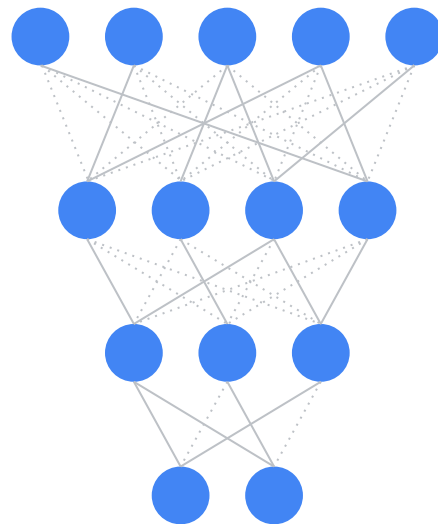
Pruning



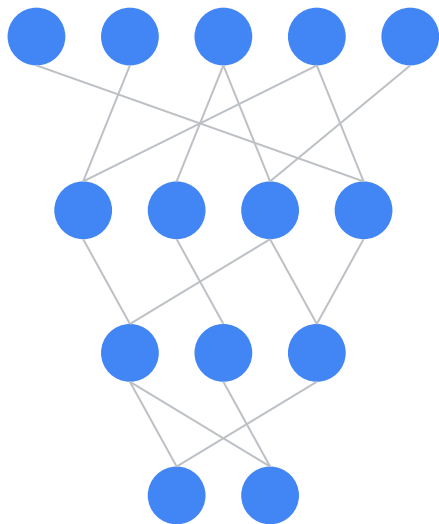
Pruning



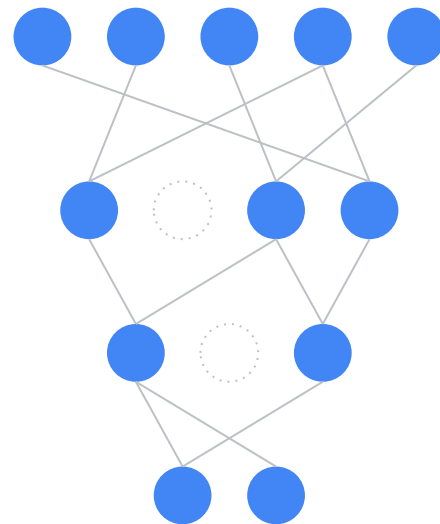
**PRUNING
SYNAPSES**



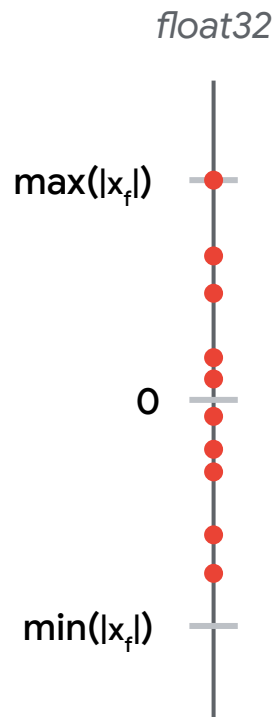
Pruning



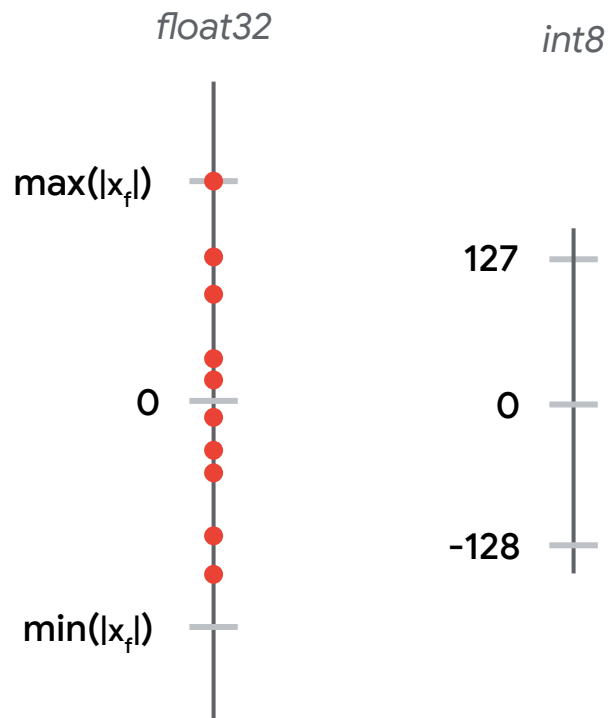
PRUNING
NEURONS



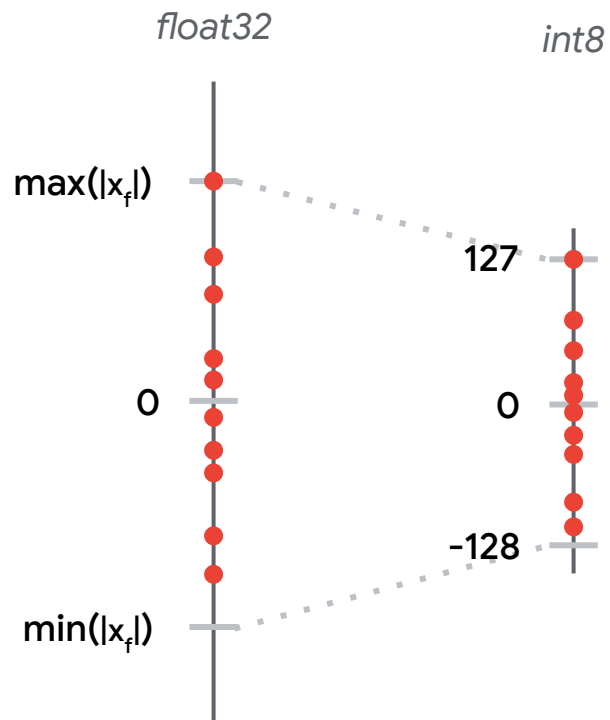
Quantization



Quantization



Quantization



1