{ Model size

↓

number of FLOPs

| ↑ Compression ratio

Maintain accuracy

**Compact Model:** Modify the standard operations used in DNNs.

(num of parameters)

e.g.

Standard CNN ← ↓ { dilated conv

separable depthwise conv

Standard LSTM ← { S-LSTM

JANET

**Tensor Decomposition:**

$$M = AB$$

$M \times n \qquad M \times r \quad r \times n$

$(M \times r) + (r \times n) < M \times n$

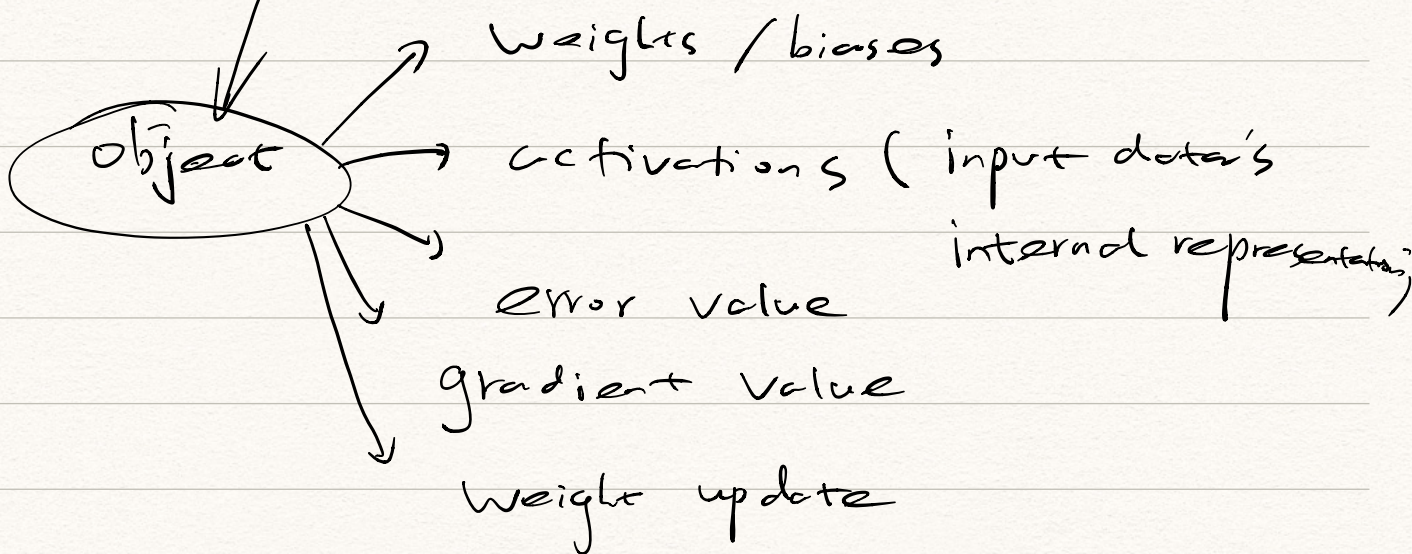hierarchical tensor representation
(HT)

Tensor train decomposition
(TT)

**Quantization :** convert ~~data~~ (object)
from 32 - floating point to
lower precision or a fixed
point integer or even
binary

Object → Weights / biases
→ activations ( input data's
          internal representation)
→ error value
→ gradient value
→ weight update

**Network Sparsification / Pruning :**

compress the model by pruning some
weights (edges) or operations (nodes)

" importance" ← Weight values

← Learned via an
Attention Layer