

Energy-Aware Neural Architecture Optimization with Splitting Steepest Descent

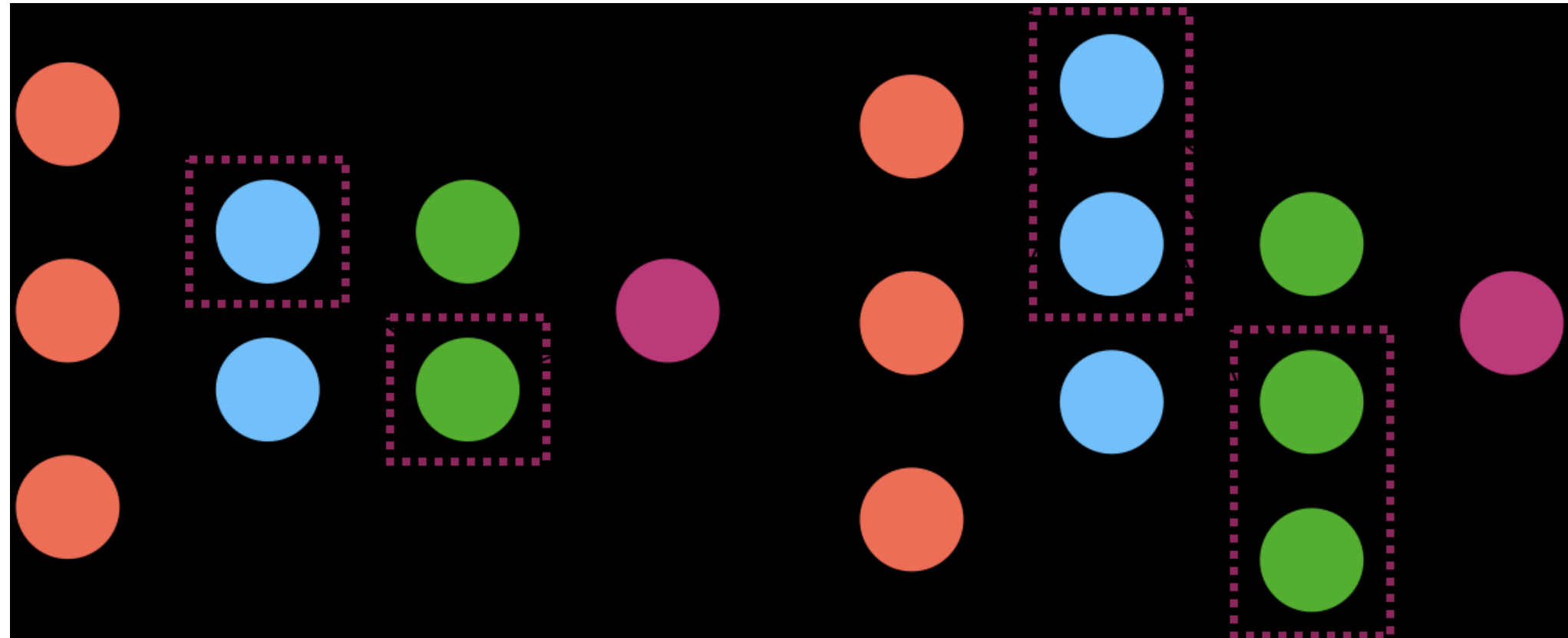
Dilin Wang¹, Lemeng Wu¹, Meng Li², Vikas Chandra², Qiang Liu¹
¹ UT Austin ² Facebook



Neural architecture optimization

Splitting Yields Adaptive Net Structure Optimization

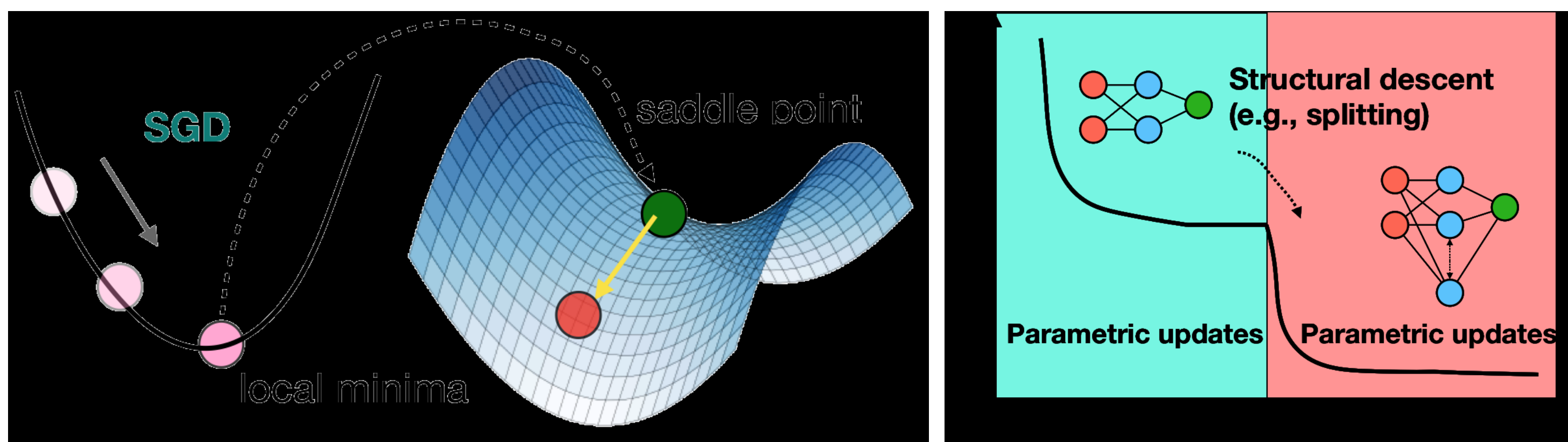
- Starting from a small net, gradually grow the net during training.
- Grow by "splitting" existing neurons into multiple offspring.



Why, when and how?

- Why splitting? Does splitting decrease the loss? How much?
- When to split? What neurons should be split first?
- How to split a neuron optimally? How many copies to split into?

Why & when: escaping local minima



- Optimization view:** the local optima in the low dimensional space can be turned into a saddle point in a higher dimensional of the augmented networks
- Architecture view:** lower-dimensional space / smaller networks; higher-dimensional space / larger networks

How: splitting steepest descent

- Consider a single-neuron network

$$L(\theta) := E_{x \sim D} [\ell(\theta; x)];$$

where $\ell(\cdot)$ is the map from the output of the neuron to the neural loss.

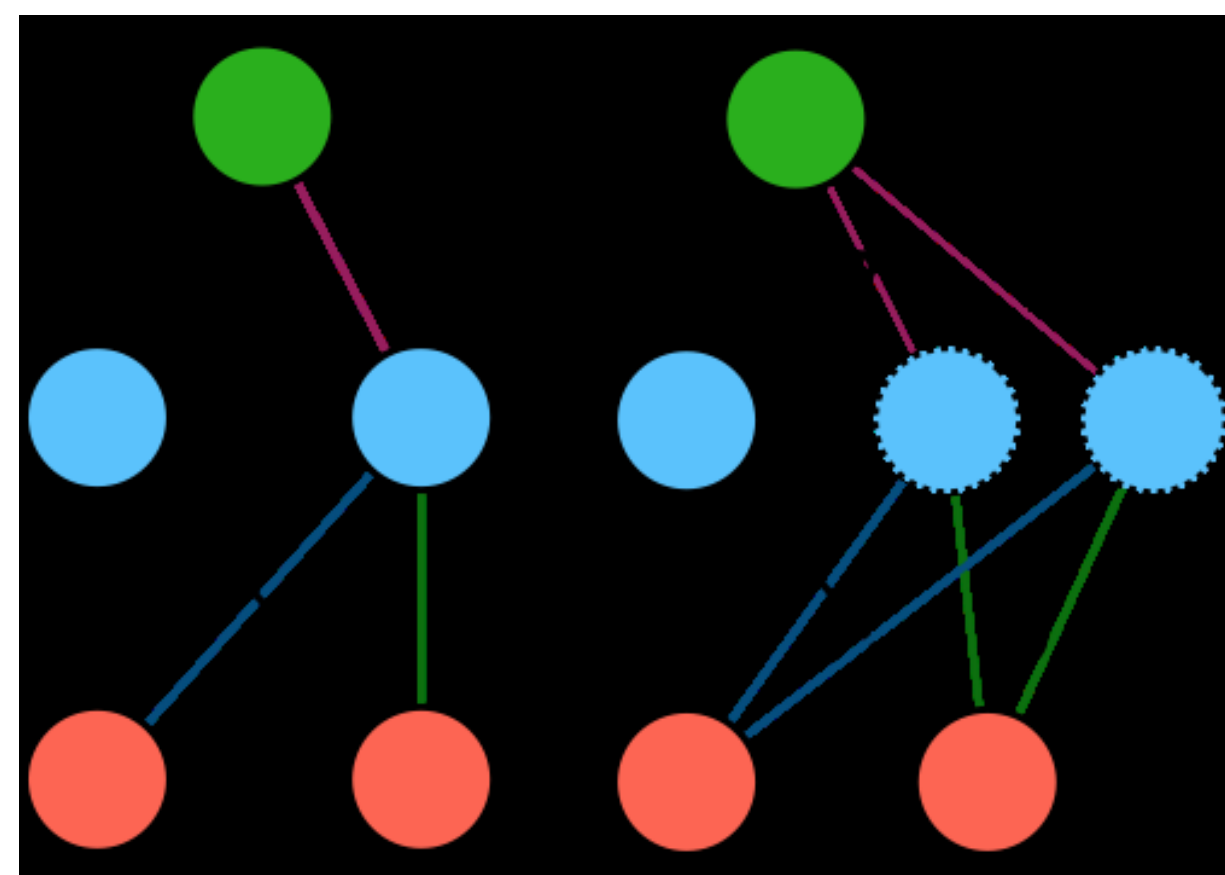
- Split into m offspring:

! $f_i; w_i g_{i=1}^m$, we have,

$$L(f_i; w_i g) := E_{x \sim D} \left(\sum_{i=1}^m w_i \ell(f_i; x) \right);$$

- Smooth loss change:**

$$w_i = 1; \quad j \neq i \quad j \neq 2 \quad ; \delta_i;$$



Deriving optimal splitting strategies

- Structural descent at stable local minima

$$\min_{m, f_i g; w_i g} L_m(f_i; w_i g) \quad L(\cdot); s.t: j \neq i \quad j \neq 2 \quad ; \quad w_i = 1; w_i > 0; \quad (1)$$

- The optimum of Eqn. 1 is determined by

$$\min_{m, f_i g; w_i g} L_m(f_i; w_i g) \quad L(\cdot) = \frac{2}{2} \min_{\text{splitting index}} \left(\frac{S(\cdot)}{Z} \right); 0g + O(\epsilon^3);$$

$$\text{with } S(\cdot) = E_{x \sim D} \left(\frac{r(\cdot; x)}{r^2(\cdot; x)} \right);$$

where $\min(S(\cdot))$ denotes the minimum eigenvalue of $S(\cdot)$.

Optimal splitting

- When $\min(S(\cdot)) = 0$, no splitting

- When $\min(S(\cdot)) < 0$:

$$m = 2; \quad \alpha = -\nu_{\min}(S(\cdot)); \quad \beta = \nu_{\min}(S(\cdot)); \quad w_1 = w_2 = 1=2;$$

The corresponding maximum decrease of loss is $-\frac{2}{2} \min(S(\cdot)) = -\min(S(\cdot))$.

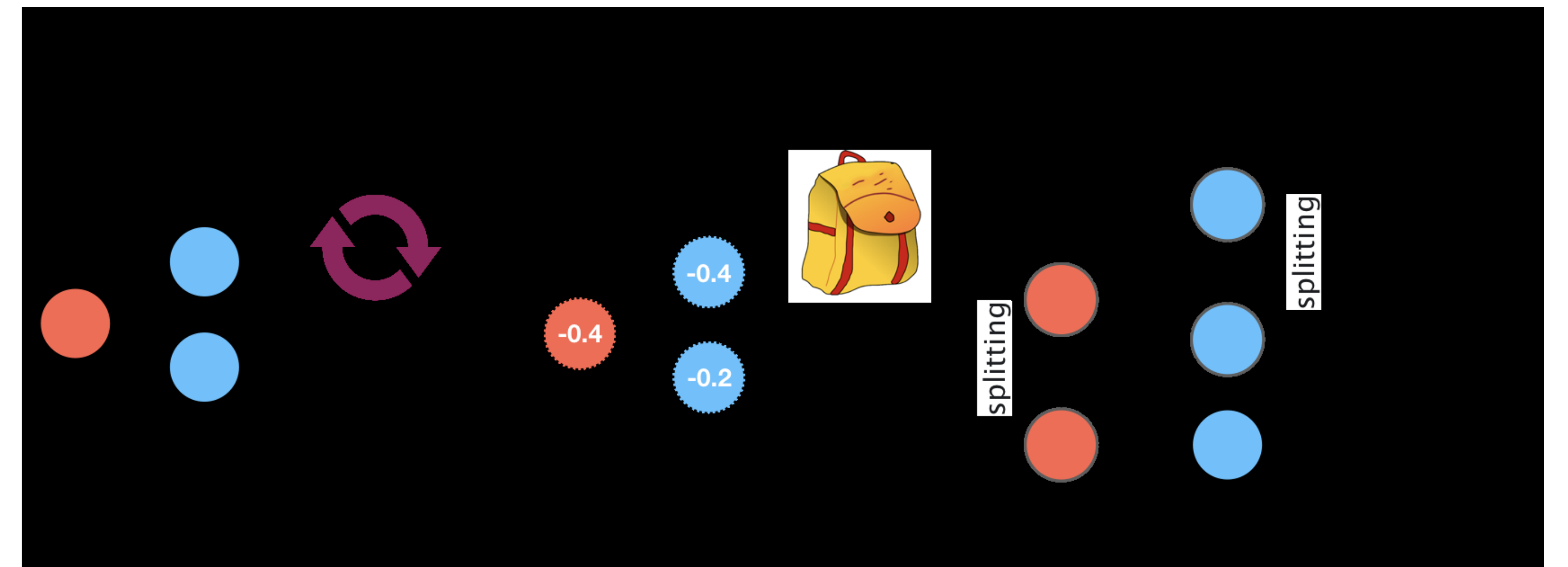
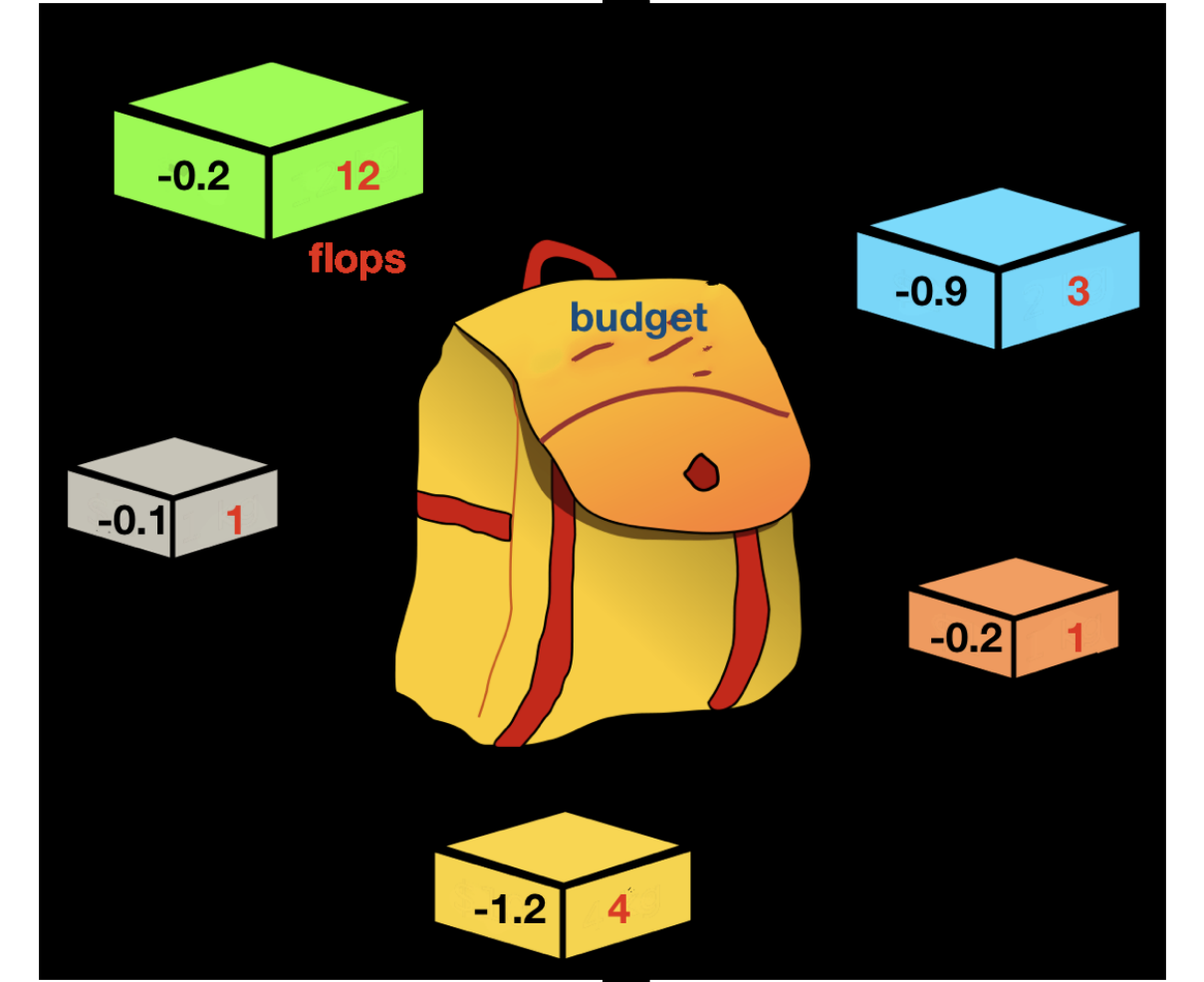
Energy-aware splitting

- Our formulation

$$\min_{\theta} \sum_{i=1}^n \left| \frac{\min\{S(\cdot)\}}{\text{gain}} \right|$$

$$s.t: \sum_{i=1}^n \left| \frac{\partial L}{\partial \theta} \right| \leq \text{budget}$$

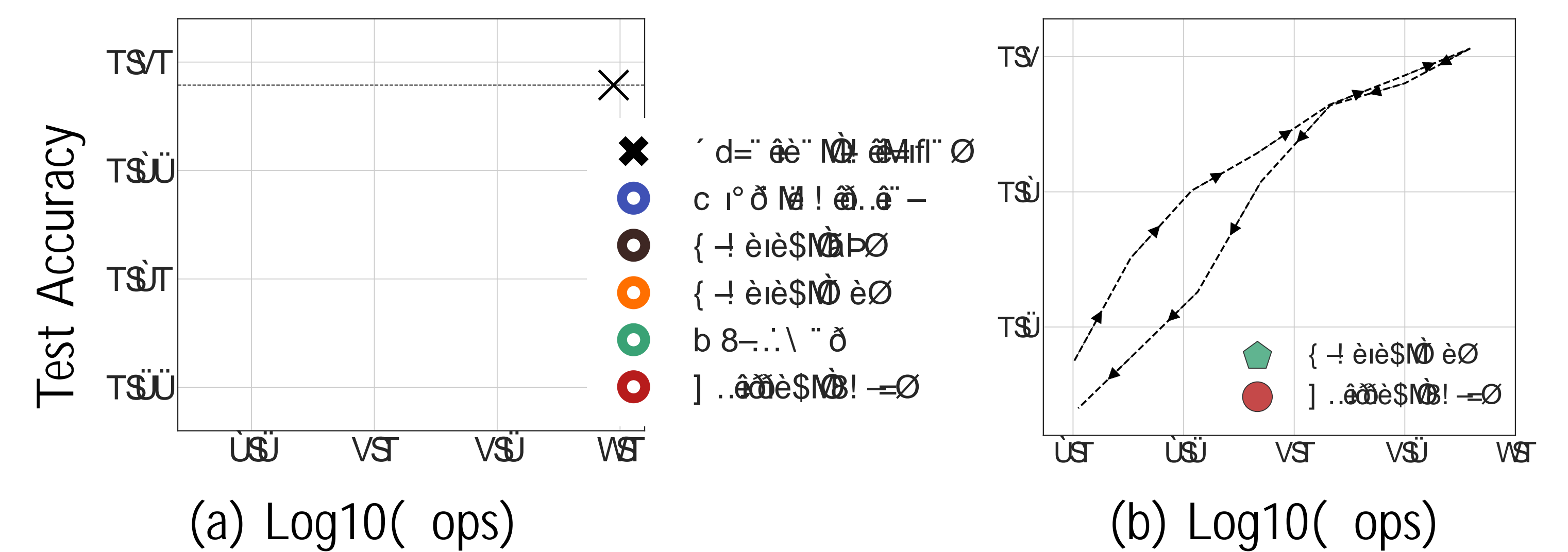
Main algorithm



Experiments

Results on CIFAR100

- We apply splitting on a small version of MobileNetV1 Howard et al., 2017



Results on ImageNet

- MobileNetV1

Model	MACs (G)	Top-1 Accuracy	Top-5 Accuracy
MobileNetV1 (1.0x)	0.569	72.93	91.14
Splitting-4	0.561	73.96	91.49
MobileNetV1 (0.75x)	0.317	70.25	89.49
AMC He et al., 2018	0.301	70.50	89.30
Splitting-3	0.292	71.47	89.67
MobileNetV1 (0.5x)	0.150	65.20	86.34
Splitting-2	0.140	68.26	87.93
Splitting-1	0.082	64.06	85.30
Splitting-0 (seed)	0.059	59.20	81.82

- MobileNetV2

Model	MACs (G)	Top-1 Accuracy	Top-5 Accuracy
MobileNetV2 (1.0x)	0.300	72.04	90.57
Splitting-3	0.298	72.84	90.83
MobileNetV2 (0.75x)	0.209	69.80	89.60
AMC He et al., 2018	0.210	70.85	89.91
Splitting-2	0.208	71.76	90.07
MobileNetV2 (0.5x)	0.097	65.40	86.40
Splitting-1	0.095	66.53	87.00
Splitting-0 (seed)	0.039	55.61	79.55

Conclusion

- Incremental training with splitting gradient.
- Simple and fast, promising in practice.
- Opens a new dimension for energy-efficient NAS.

