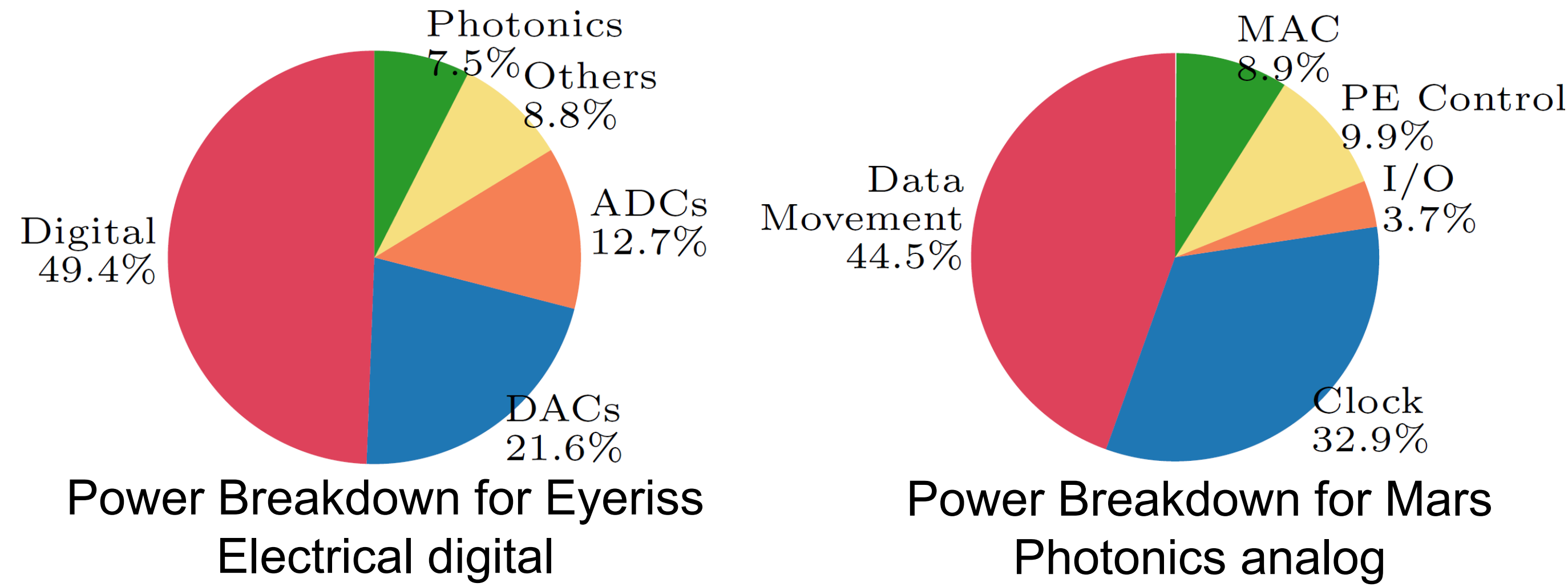
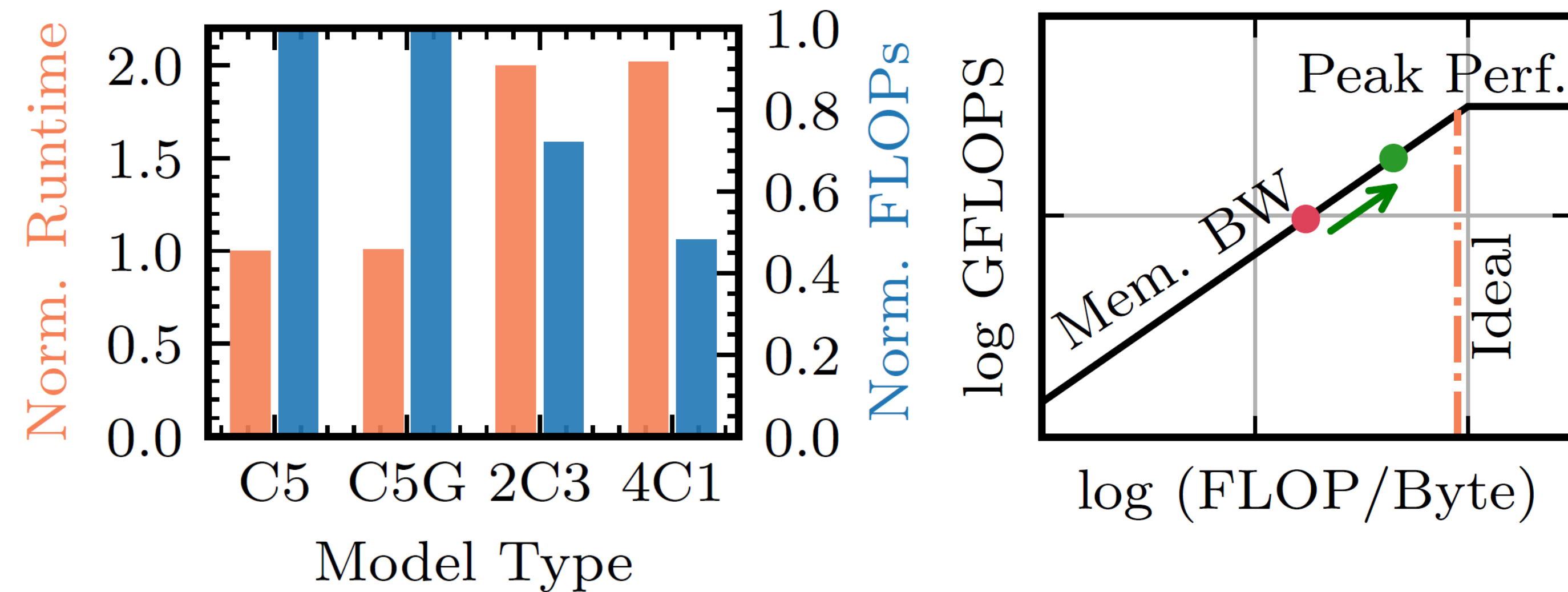


## Introduction:

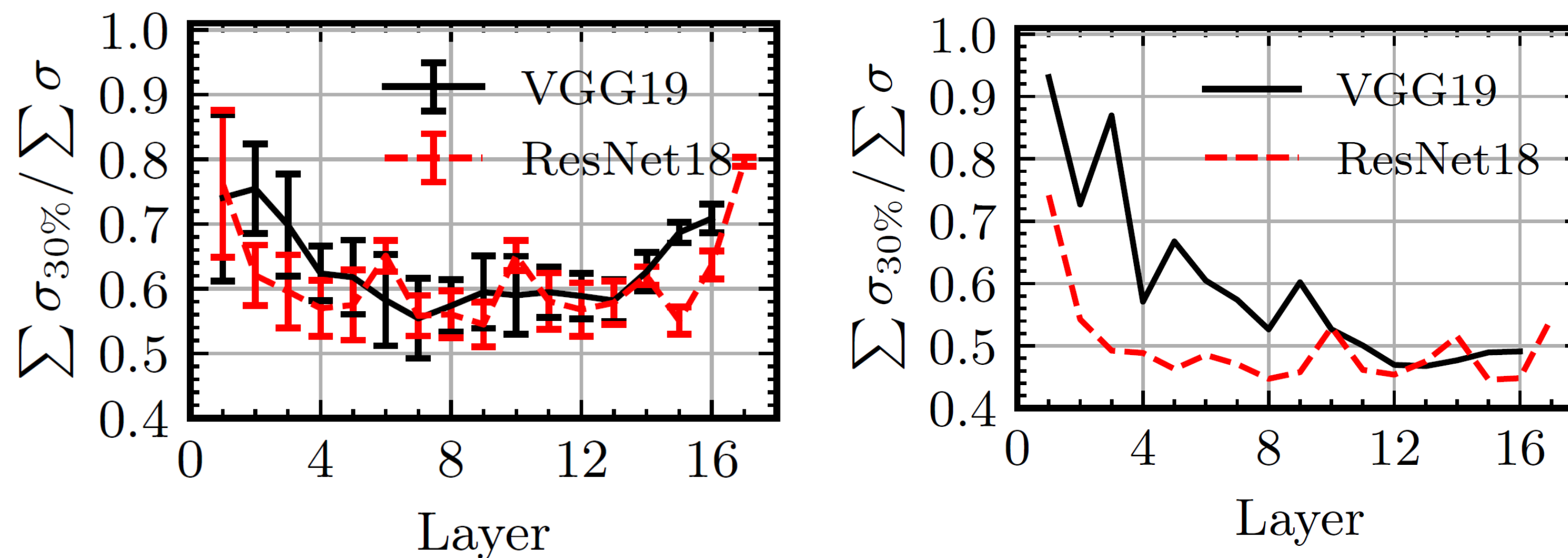
- Data movement is more costly than computations



- Memory bottlenecked modern AI accelerators efficiency
- Naive layer decomposition does not save memory cost

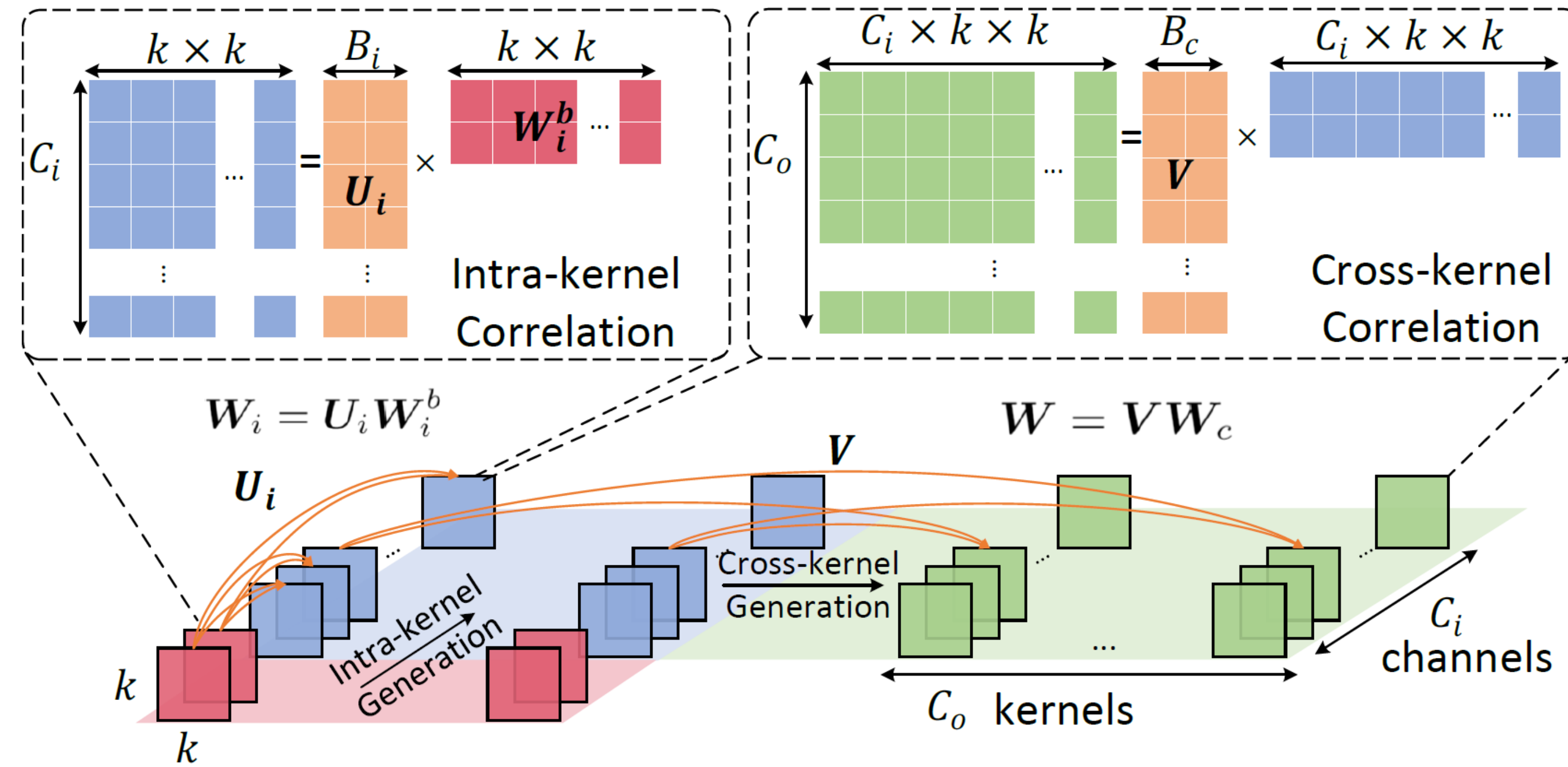


- Intrinsic intra-kernel and cross-kernel correlation in CNNs
- Motivate us to generate kernel on-the-fly with small basis



## Proposed Multi-Level *in-situ* Generation:

- Intra-kernel generation  $W_i = U_i W_i^b, \forall i \in [C_o]$ 
  - Span all input channels from a small basis  $W^b$
- Cross-kernel generation  $W = V W_c = V \{U_i W_i^b\}_{i \in [B_c]}$ 
  - Span all kernels from a kernel basis
- Augmented mixed-precision generation
  - Assign different bitwidth to basis and coefficient



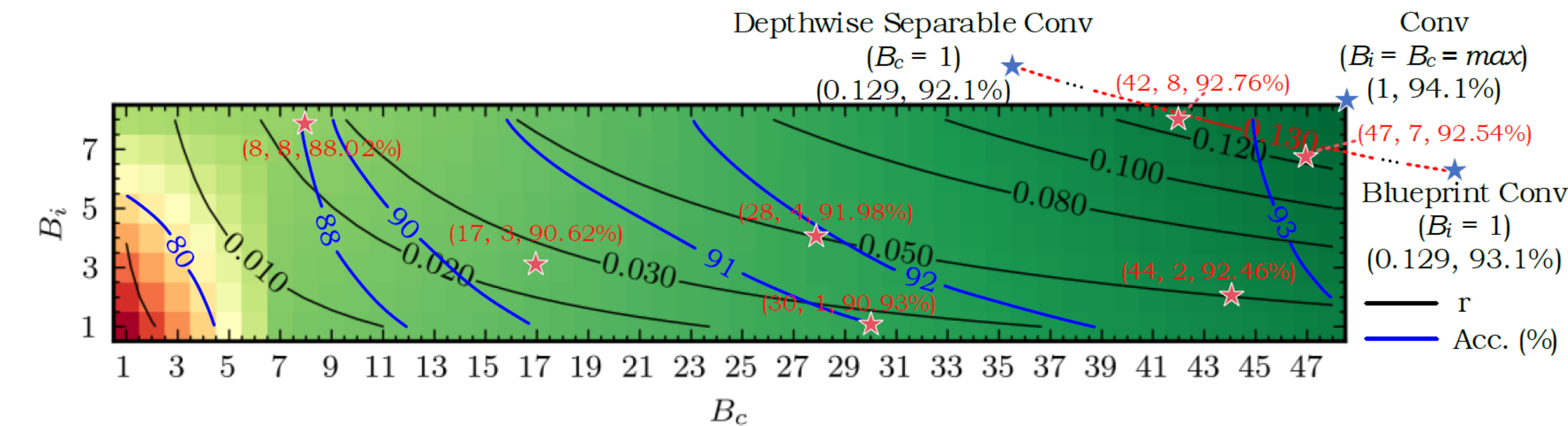
## Proposed Training Flow:

- Project teacher onto decomposed low-bit students
 
$$\min \|\widehat{\mathcal{M}}(\widehat{W}) - \mathcal{M}(W)\|_2^2 \approx \|\widehat{W} - V \{U_i W_i^b\}_{i \in [B_c]}\|_2^2$$
- Distill knowledge from teacher to students
 
$$\min \mathcal{L}_{KD} = \beta T^2 \mathcal{D}_{KL}(q_T, p_T) + (1 - \beta) H(q, p_{T=1})$$

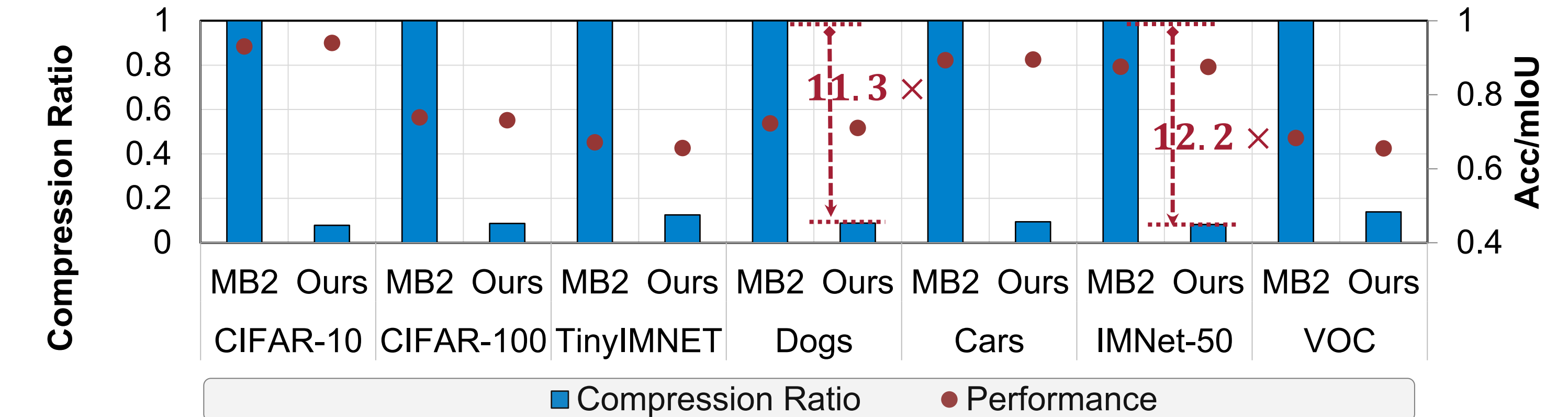
$$\text{s.t. } p_T = \frac{\exp(\frac{\mathcal{M}(W)}{T})}{\sum \exp(\frac{\mathcal{M}(W)}{T})}, q_T = \frac{\exp(\frac{\widehat{\mathcal{M}}(\widehat{W})}{T})}{\sum \exp(\frac{\widehat{\mathcal{M}}(\widehat{W})}{T})}$$
- Orthonormal regularization to encourage ranks
 
$$\sum_{i=1}^{B_c} \left( \|W_i^b (W_i^b)^T - I\|_2^2 + \|\tilde{U}_i^T \tilde{U} - I\|_2^2 \right) + \|\tilde{V}^T \tilde{V} - I\|_2^2$$

## Experimental Results:

- Intra-kernel correlation is stronger than cross-kernel correlation
- Outperforms separable CONV and Blueprint CONV



- Comparable performance on various tasks with compact models
- **>10x** memory compression (3~5 bit quantization)
- **~30%** less latency, **~86%** less energy on simulated ReRAM Accel.



## Summary/Conclusion

- Multi-level in-situ generation for memory-efficient DNN design
- Mixed-precision for fine-grained design space exploration
- **10-20x** memory compression; **~97%** less weight loading time
- New design paradigm to break through the ultimate memory bottleneck for emerging DNN accelerators