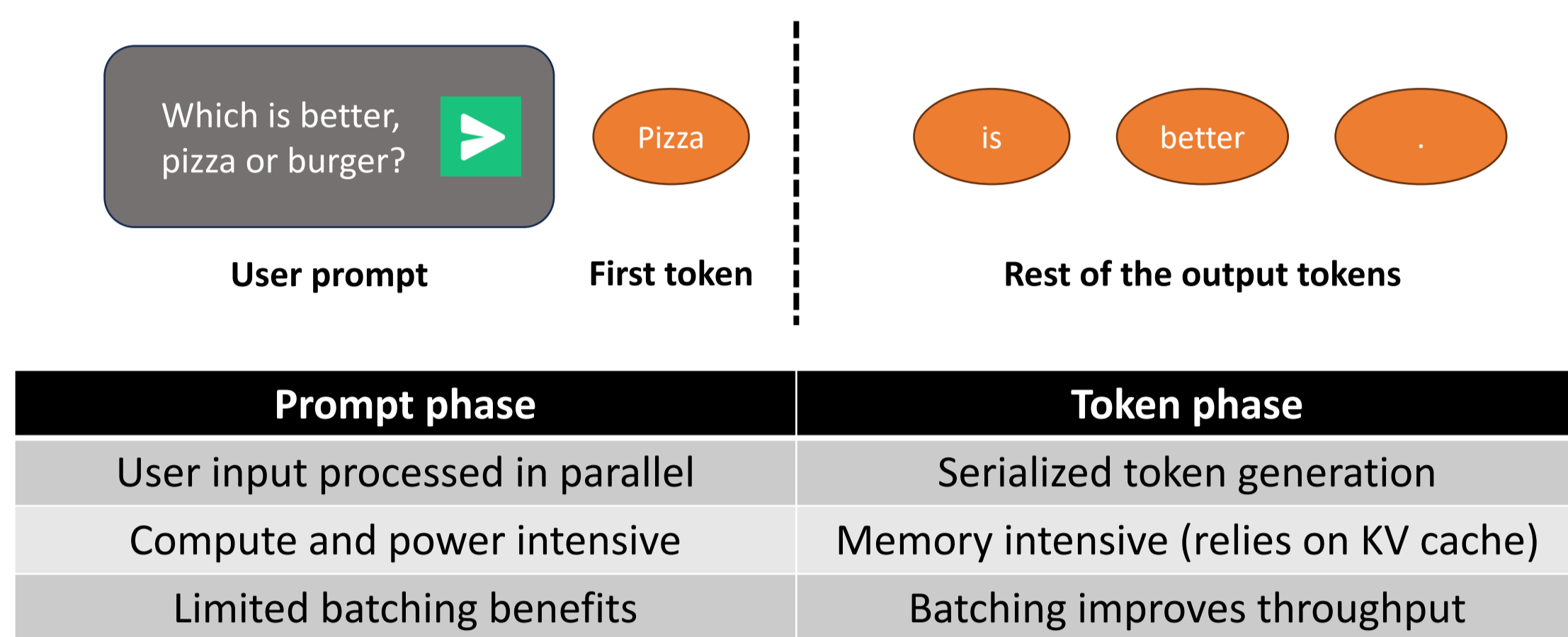


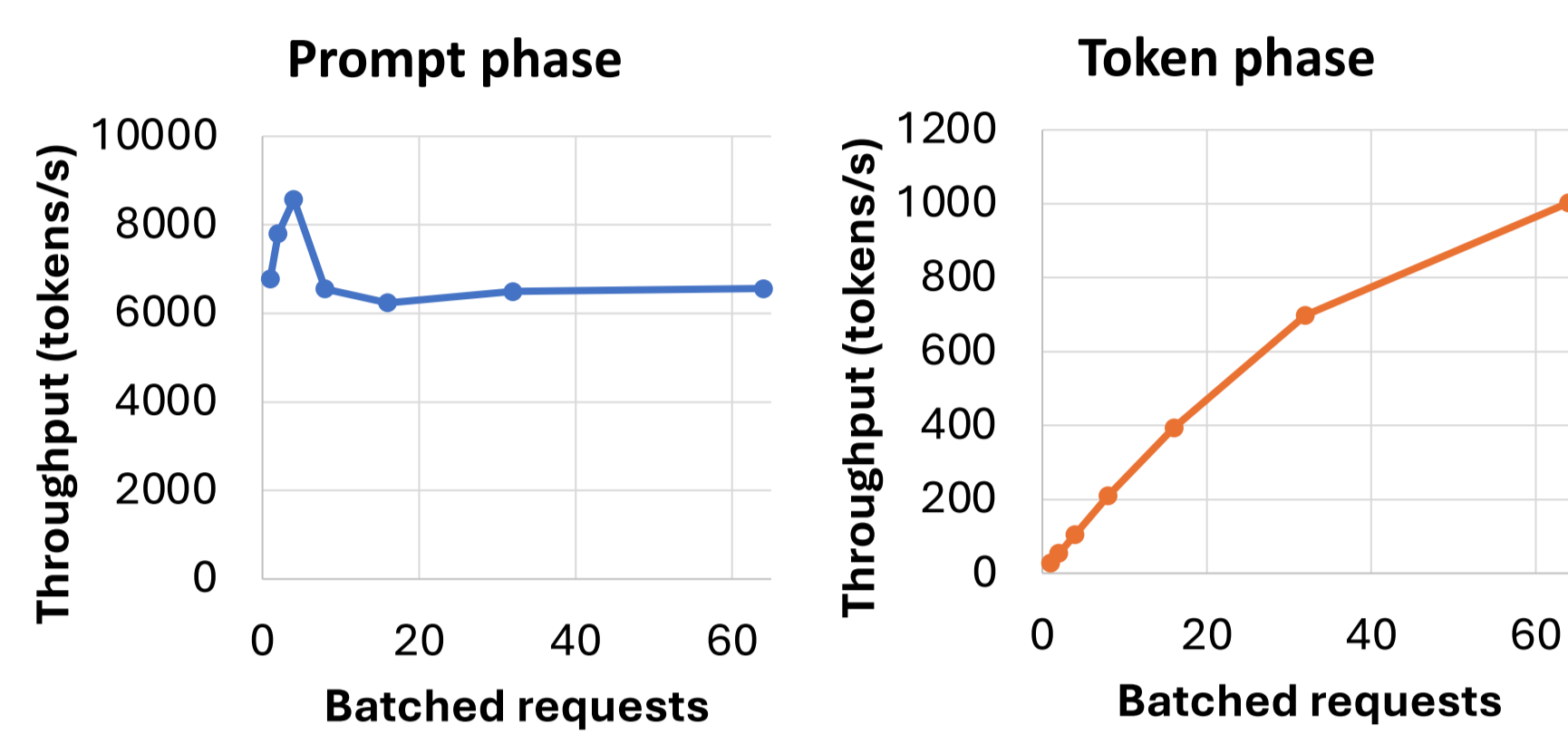
# Splitwise: Efficient Generative LLM Inference Using Phase Splitting

Each LLM inference request has two distinct phases with different resource requirements

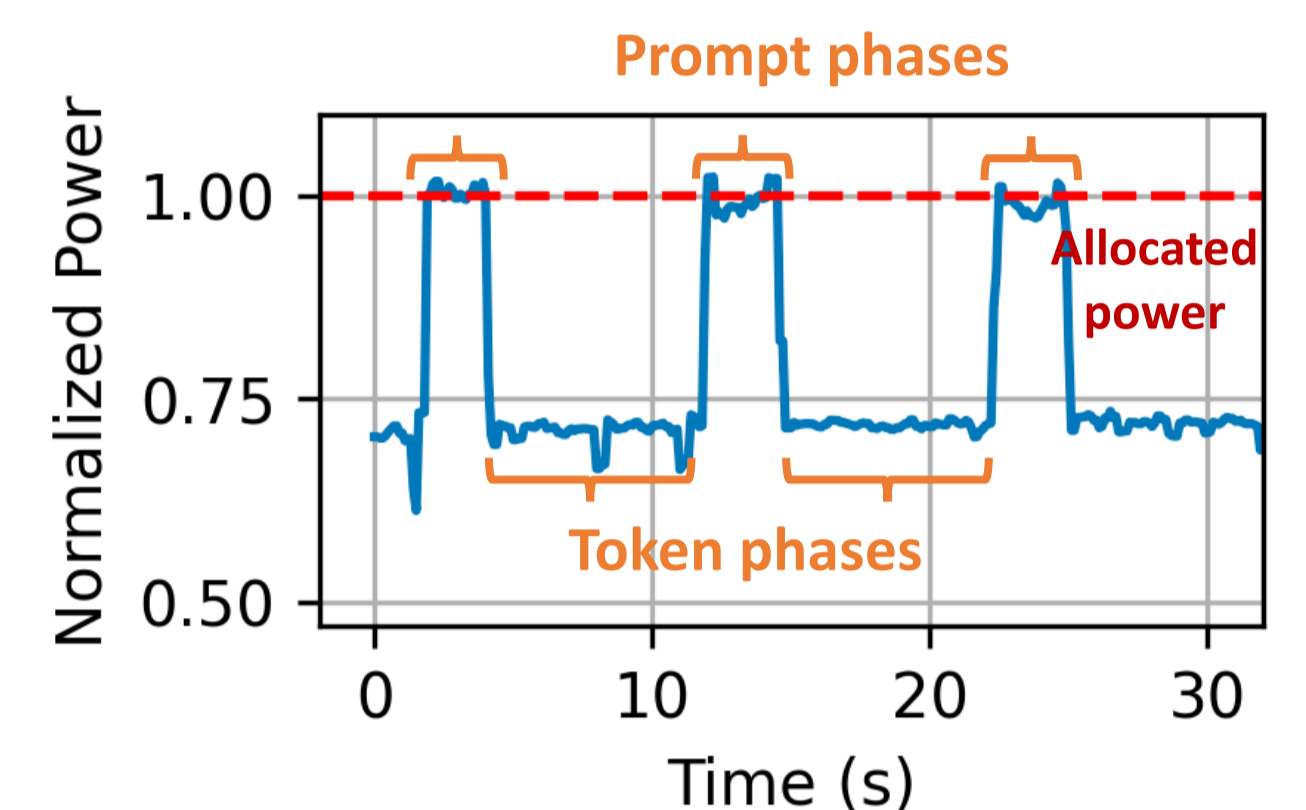
Prompt computation and token generation phases



Example 1: Batching effects

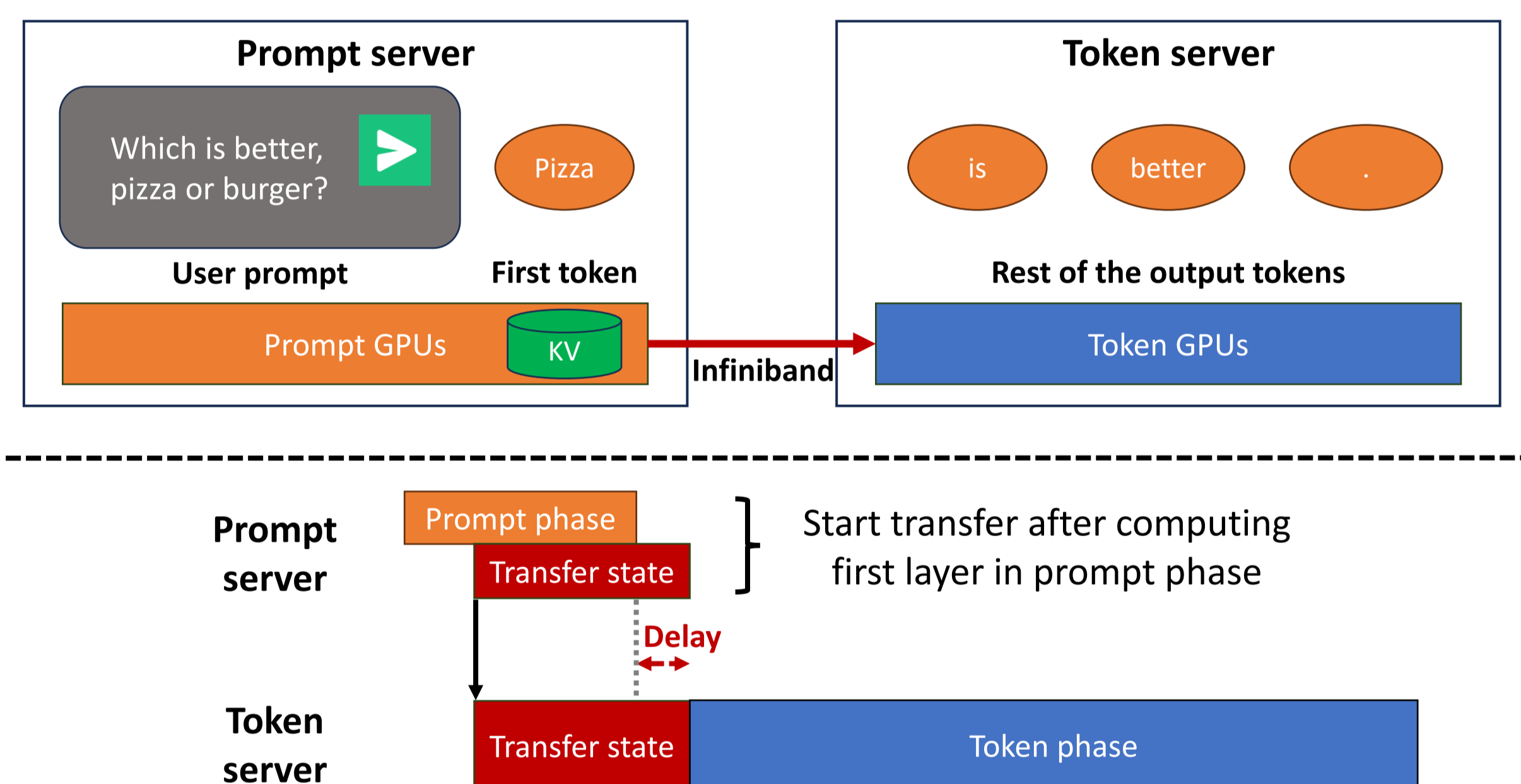


Example 2: Power usage

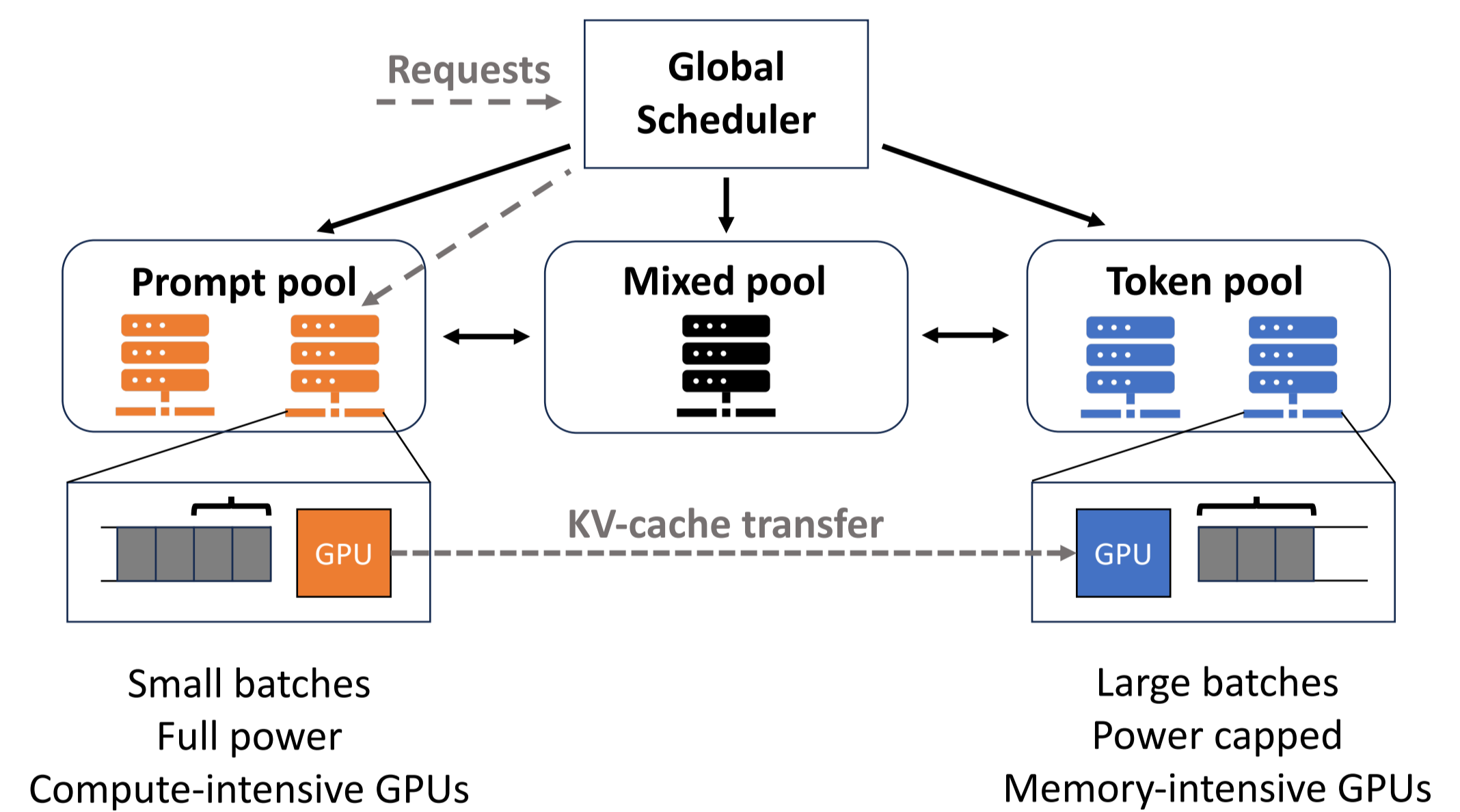


Splitwise splits inference across different servers to enable phase-specific resource management

1 Transfer request state over P2P GPU Infiniband; optimize with parallel and overlapped transfers

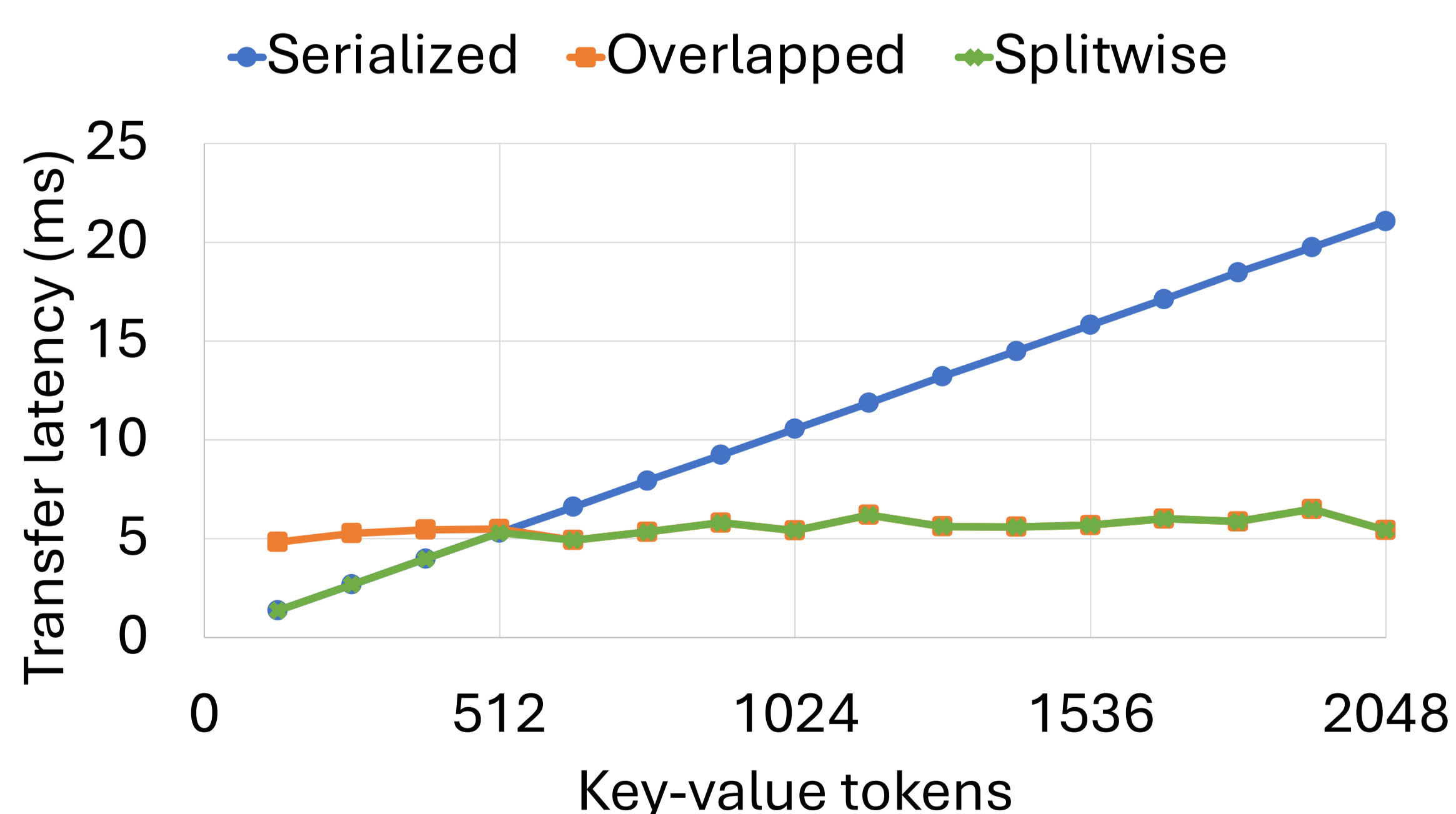


2 Split cluster into three server pools and use phase-specific resource management at scale



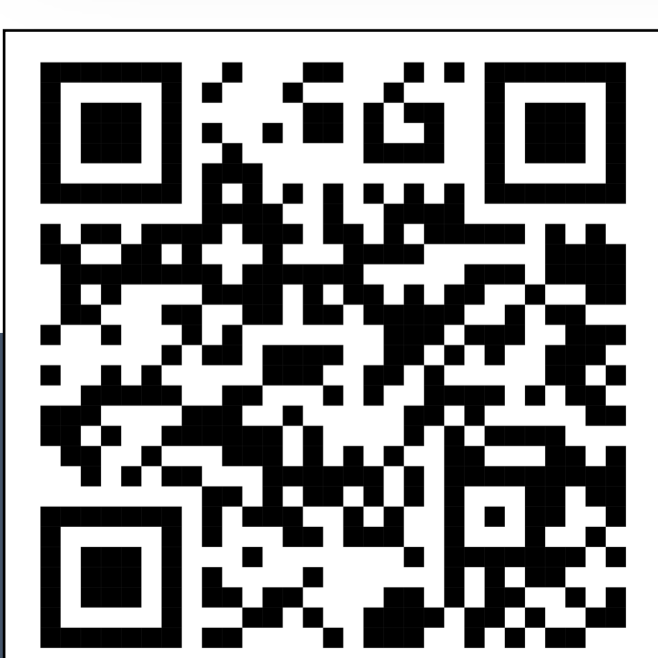
Splitwise clusters are much more resource efficient than existing clusters

Result 1: Splitwise transfers request state with less than ~0.8% end-to-end overhead on average



Result 2: Clusters designed using Splitwise provide much higher throughput than existing clusters

	Baseline	Splitwise homogeneous	Splitwise heterogeneous
Throughput optimized clusters	A100 x70	A100 x45 → A100 x25	H100 x25 → A100 x26
#Servers	1x	1x	0.73x
Cost	1x	1x	1.14x
Power	1x	1x	1x
Throughput	1x	2.4x	2.6x



Paper, code, traces at [aka.ms/splitwise](https://aka.ms/splitwise)

Pratyush Patel, Esha Choukse, Chaojie Zhang, Aashaka Shah, Ñigo Goiri, Saeed Maleki, Ricardo Bianchini

