

Parameter servers for machine learning do not scale?

Well, one size does not fit all.

NuPS: A Parameter Server for Machine Learning with Non-Uniform Parameter Access

Alexander Renz-Wieland, Rainer Gemulla, Zoi Kaoudi, Volker Markl

To appear in SIGMOD 2022

Motivation

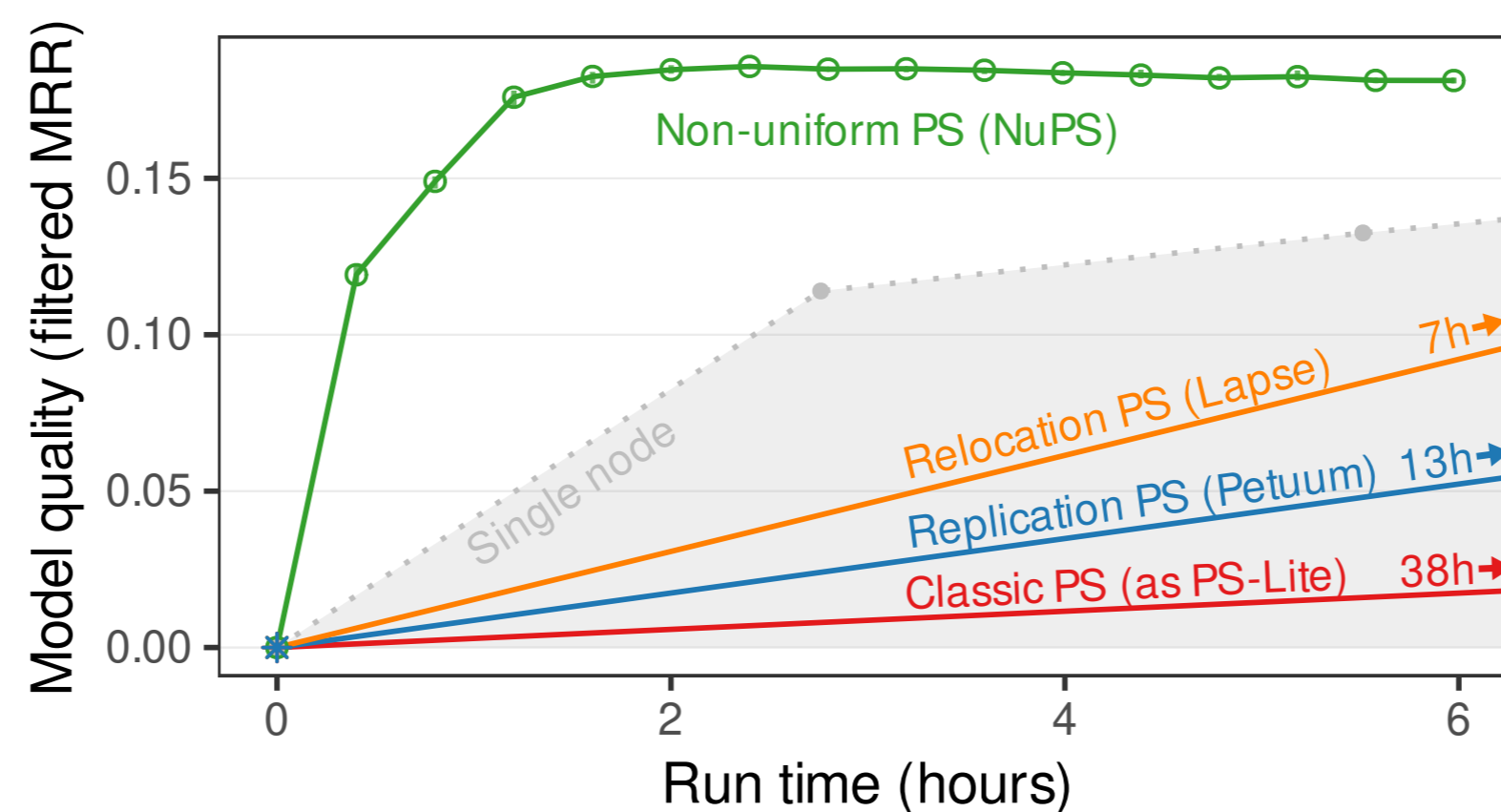
- Parameter Servers (PSs) facilitate distributed machine learning
- But existing PSs are inefficient for non-uniform access:
 - Skew: PSs inefficient because they manage all parameters identically
 - Sampling: PSs inefficient because sampling entails randomized access

Results

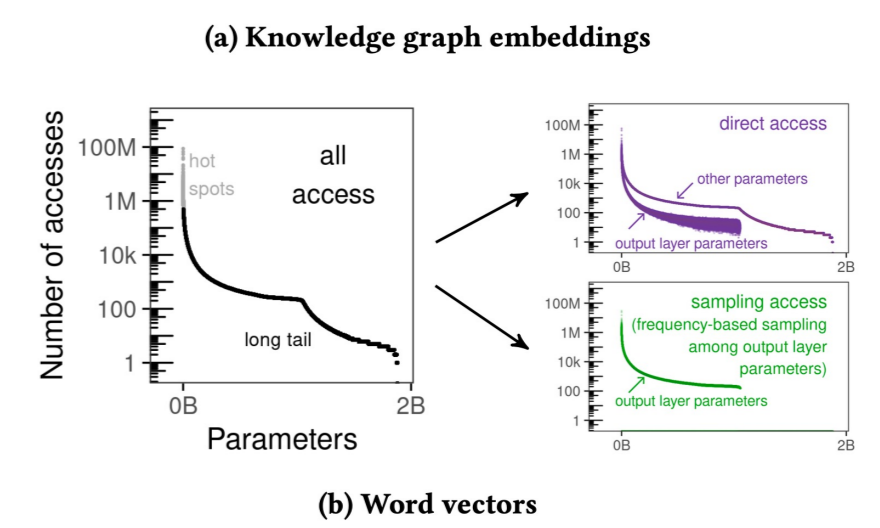
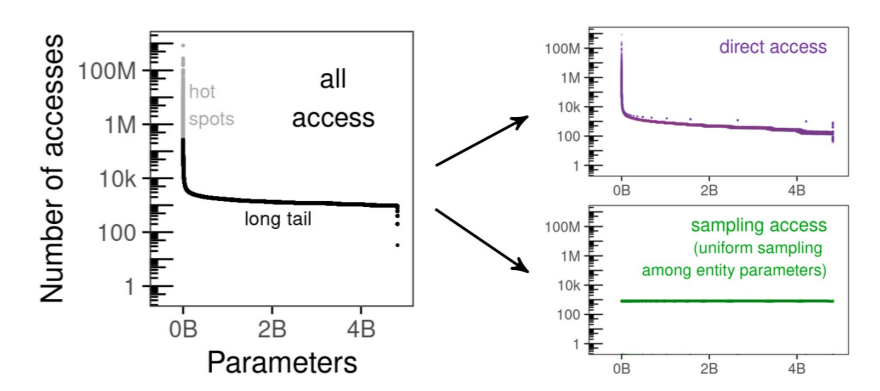
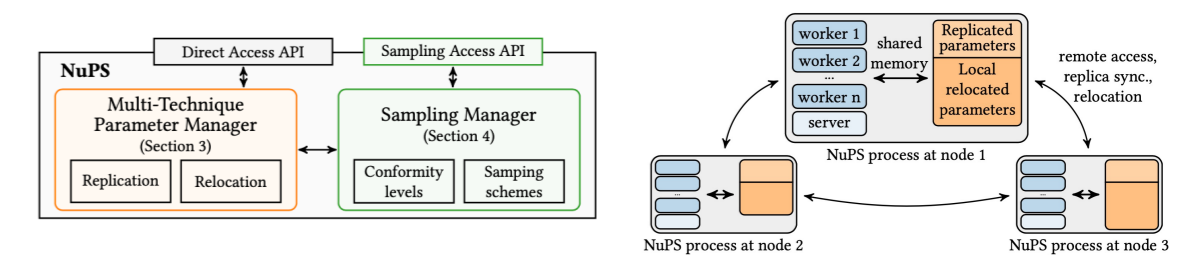
- NuPS outperformed existing PSs by up to one order of magnitude across multiple ML tasks
- NuPS provided up to linear scalability
- Most efficient was to replicate a small set of hotspot parameters and relocate all others

NuPS

- Supports multiple management techniques and picks a suitable one for each parameter:
 - Replication is efficient for hot spot parameters
 - Relocation is efficient for long tail parameters
- Supports sampling directly via suitable primitives and sampling schemes that allow for a controlled quality–efficiency trade-off



Parameter server performance for training large knowledge graph embeddings on an 8-node cluster.



handle = PrepareSample(dist, N)
keys, values = PullSample(handle, n_j)

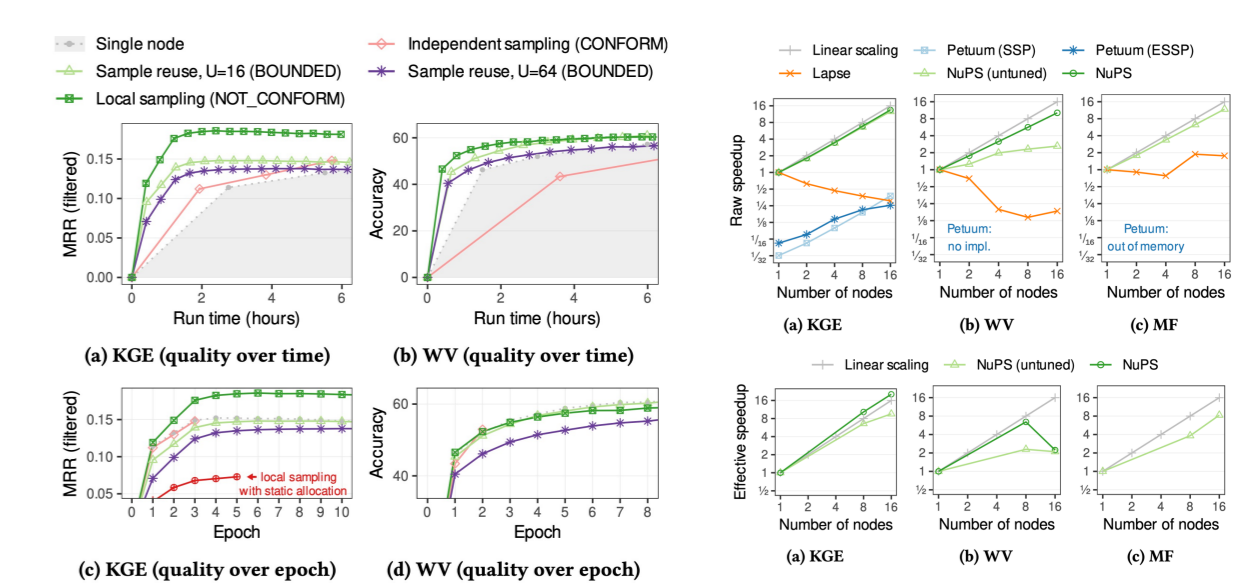
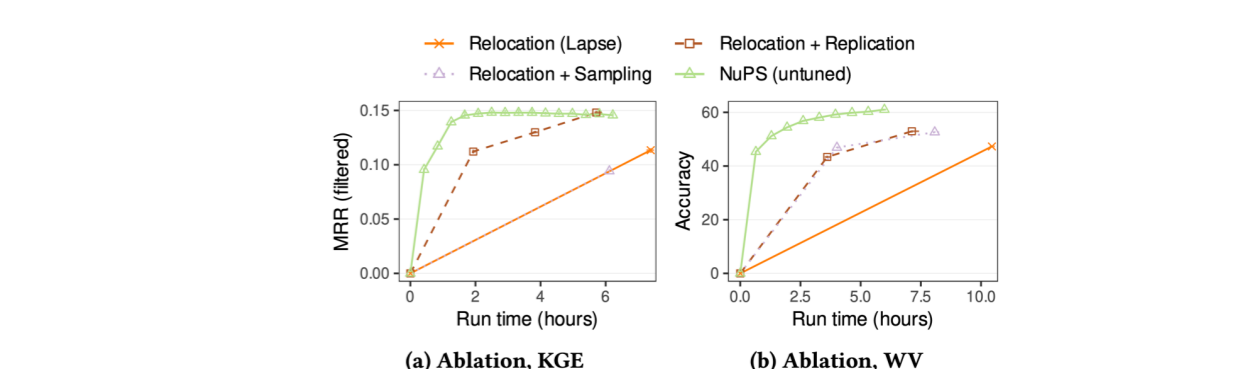
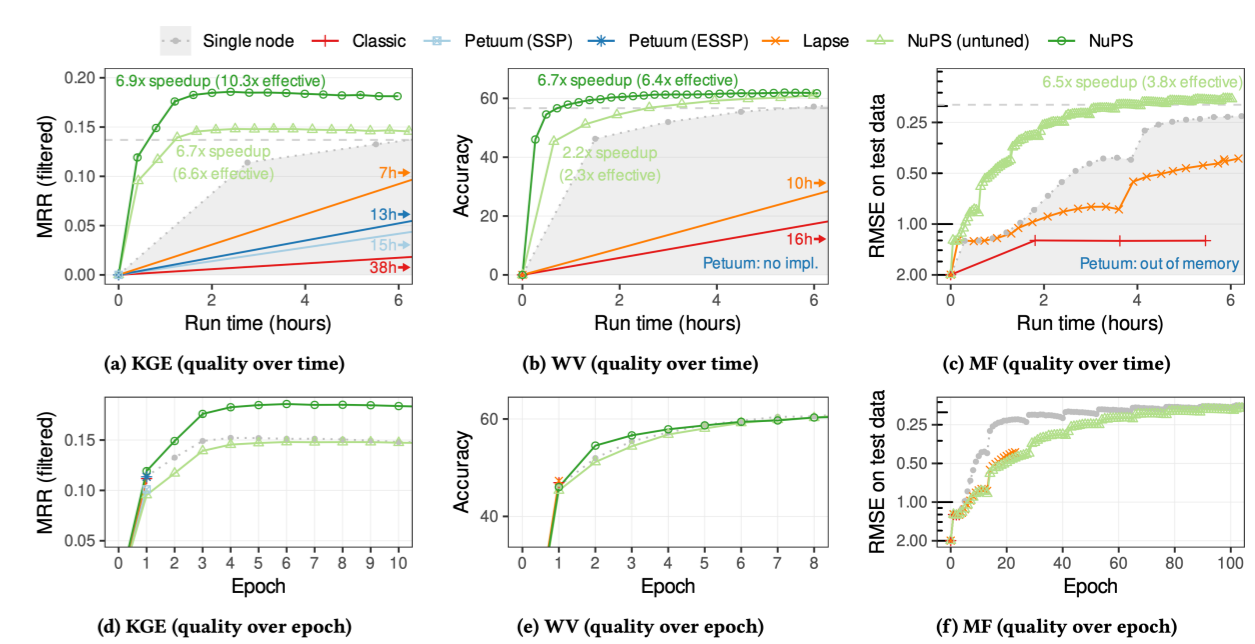
	L1	L2	L3
Independent sampling	✓	✓	✓
Sample reuse	✗	✗	✗
Local sampling	✗	✗	✗
Direct access (reputing)	✗	✗	✗

(0) CONFORM: The sampling scheme produces mutually independent samples from the target distribution π .
 (1) BOUNDED: The samples at each node have dependencies on past samples, but these dependencies are limited and samples at different nodes are independent. In more detail, given a dependency bound $\beta \in \mathbb{N}$,

$$p_{i,j}^{(k)} = \mathbb{P}(x_i^{(k)} = v_j | x_{i-\beta}^{(k)} = v_{j'}) = \pi_j$$
 for all i, j, k , where $\beta, \gamma \in \mathbb{N}$. $x_i^{(k)}$ refers to samples at other nodes and $x_{i-\beta}^{(k)}$ refers to samples at node i .
 (2) NON-CONFORM: The samples at each node have dependencies on past samples, but these dependencies are limited and samples at different nodes are independent. In more detail, given a dependency bound $\beta \in \mathbb{N}$,

$$p_{i,j}^{(k)} = \mathbb{P}(x_i^{(k)} = v_j | x_{i-\beta}^{(k)} = v_{j'}) = \pi_j$$
 for all i, j, k , where $\beta, \gamma \in \mathbb{N}$. $x_i^{(k)}$ refers to samples at other nodes and $x_{i-\beta}^{(k)}$ refers to samples at node i .
 (3) LOCAL SAMPLING: The samples at each node have dependencies on past samples, but these dependencies are limited and samples at different nodes are independent. In more detail, given a dependency bound $\beta \in \mathbb{N}$,

$$p_{i,j}^{(k)} = \mathbb{P}(x_i^{(k)} = v_j | x_{i-\beta}^{(k)} = v_{j'}) = \pi_j$$
 for all i, j, k , where $\beta, \gamma \in \mathbb{N}$. $x_i^{(k)}$ refers to samples at other nodes and $x_{i-\beta}^{(k)}$ refers to samples at node i .
 (4) NON-CONFORM: No guarantees about the sampling probabilities at independence.



Task	Model parameters			Data		Parameter access	
	Model	Keys	Values	Size	Dataset	Data points	Size
Knowledge graph embeddings	ComplEx, dim. 500	4.8 M	4.8 B	35.9 GB	Wikidata5M	21 M	317 MB
Word vectors	Word2Vec, dim. 1000	1.9 M	1.9 B	7.0 GB	1b word benchmark	375 M	3 GB
Matrix factorization	Latent Factors, rank 1000	11.0 M	11 B	82.0 GB	10m x 1m matrix, zipf 1.1	1000 M	31 GB



Read the full paper

