

String Figure: A Scalable and Elastic Memory Network Architecture

Matheus Ogleari

Ye Yu, Chen Qian, Ethan Miller, Jishen Zhao

HPCA 2019





Issues with Memory Capacity Bottleneck

Training neural networks is data intensive, and increasingly so.



Why Memory Network

Memory Networks



Why Memory Network

- Memory Networks
- More sockets, more problems



A Quad-socket Server System



Memory Node

Processor

Scalability

Arbitrary Network Scale

Elastic Network Scale

















String Figure: A Scalable and Elastic Memory Network Design

Our Goals

Scalability

Arbitrary Network Scale

Elastic Network Scale



Topology design



Random topology generation



Short cuts

Scalability

Arbitrary Network Scale

Elastic Network Scale





Greedy routing protocol

Node #	<i>D</i> in Space-0	MD	
7	0.49	0.62	0.49
0	0.20	0.70	0.20
3	0.13	0.44	0.13
6	0.43	0.12	0.12
8	0.68	0.07	0.07

Node#	Block	. Valid	Нор	Space#	Coordi.
0	1	1	0	0	0.00
2	1	1	0	0	0.20
5	1	1	0	1	0.58
6	1	1	0	1	0.75
				•	
3	1	1	1	0	0.33
8	1	1	1	0	0.88

Minimum distance (MD) to Node-2 from Node-7 and Node-7's neighbors

Routing table entries

Implementation



Design Parameters:

- Payload Size (bits)
- # of net ports
- # of routers
- # of router ports

Simulation



Evaluation

	Number of Nodes (N), Number of Ports per Router (p)												
Topology	Ν	16	17	32	61	64	113	128	256	512	1024	1296	Routing Scheme
Distributed-Mesh (DM) /													
Optimized DM (ODM)	р	4	n/a	4	n/a	4	n/a	4	4	4	4	4	Greedy + adaptive
Flattened Butterfly (FB)	р								20	24	31	33	minimal + adaptive
Adaptive FB (AFB)	р								13	17	23	25	minimal + adaptive
Space Shuffle Ideal (S2-ideal)	р	4	4	4	4	4	4	4	8	8	8	8	look-up table (LUT)
String Figure (SF)	р	4	4	4	4	4	4	4	8	8	8	8	LUT + greedy + adaptive

Experiment Setup

- Traffic Patterns specific routing behaviors in networks
 - Uniform Random
 - Tornado
 - Hotspot
 - Opposite

- Nearest Neighbor
- Complement
- Partition

- Workloads real-world applications for memory networks
 - Spark-wordcount
 - Spark-grep
 - Spark-sort
 - Pagerank

- Redis
- Memcached
- Matrix Multiply
- Kmeans











- More in the paper!
 - Traffic pattern latencies
 - Real workload
 performance
 - Network saturation
 - Energy-Delay Product (EDP)
 - Network scaling
 - Deadlock avoidance

Summary of String Figure

Benefits

Scalability

Arbitrary Network Scale

Elastic Network Scale





String Figure: A Scalable and Elastic Memory Network Architecture

Matheus Ogleari

Ye Yu, Chen Qian, Ethan Miller, Jishen Zhao

HPCA 2019







