

# A Reconfigurable, Distributed Memory Accelerator for Sparse Applications

Courtney Golden (MIT CSAIL), Joel S. Emer (MIT CSAIL / NVIDIA), and Daniel Sanchez (MIT CSAIL) {cgolden, emer, sanchez}@csail.mit.edu

Motivation & Background		
<b>Characteristics of Scientific and Graph Workloads</b>		
Arithmetic Intensity	<ol> <li>Iterative</li> <li>Highly sparse</li> <li>Static and dynamic sparsity</li> <li>Low intra-iteration arithmetic intensity</li> </ol>	We cas Eins by dat acr the
<u>Shortcomings of Prior Work</u>		ope
Prior all-SRAM architectures overcome the memory bottleneck but suffer from <i>low programmability</i> or <i>low</i>		oca san
Einsum Notation		Ex:
output tensor rank names $Z_{m,n}^{M,N} = A_{m,k}^{M,K} \cdot B_{k,n}^{K,N}$ rank-2 tensor index into rank M generic dot operator		Z <sub>m</sub> <sup>M</sup> stric sho duplic

#### Data Partitioning

Partitioning data across distributed memories has a large performance impact because it directly determines network traffic and load balance.

> Objectives for data partitioning: 1. Load balance work among tiles 2. Minimize inter-tile communication

We partition data at a single-element granularity, allowing any nonzero of any tensor to be placed at any tile.



co-locate elements in the same **row** 



co-locate elements in the same **column** 



## **Programming Model**



- 1. Identify data that are likely to be used in the same iteration and place into clusters.
- 2. Independently partition each cluster over tiles, minimizing intra-cluster communication.





### Hardware Architecture



#### Results





