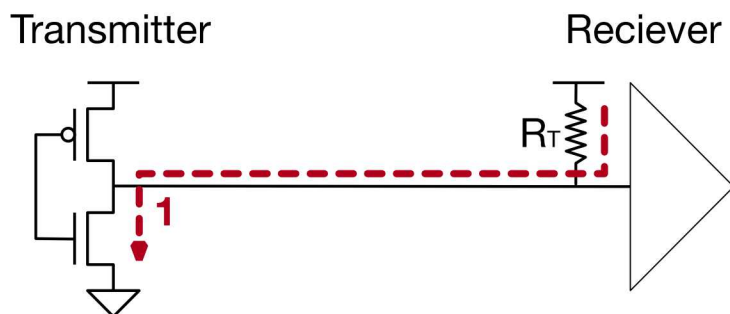# Commutative Data Reordering to Reduce Data Movement Energy on Sparse Inference Workloads
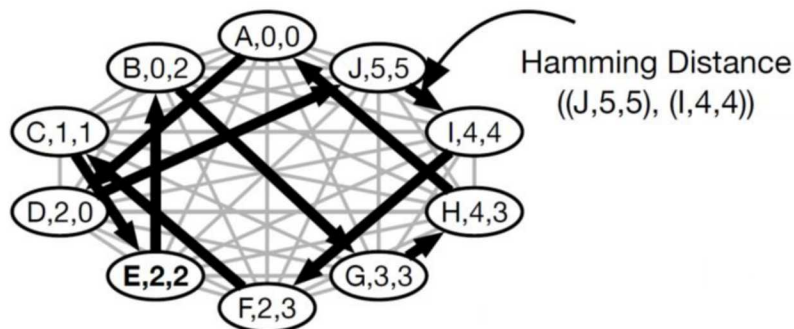
SAND2020-5473C

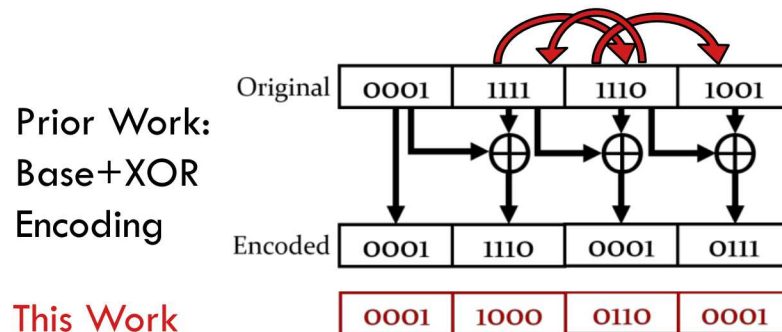Ben Feinberg[*], Benjamin C. Heyman[†], Darya Mikhailenko[†], Ryan Wong[†], An C. Ho[†], Engin Ipek[†]
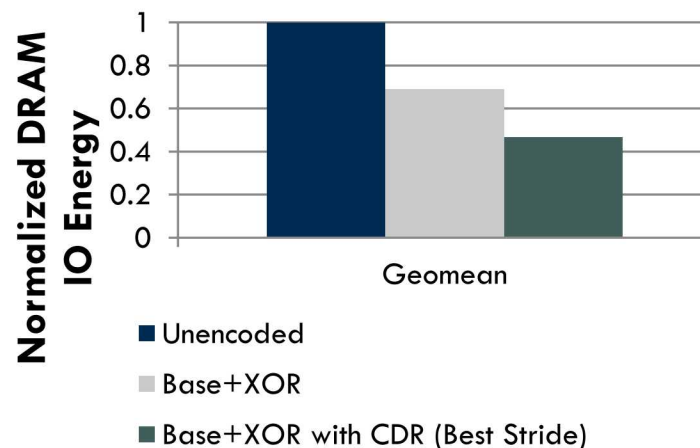
* Sandia National Laboratories, † University of Rochester

Transmitter                    Reciever



**Reducing DRAM IO energy can be modeled as reducing transmitted 1s**

Prior Work: Base+XOR Encoding

<span style="color:red">This Work</span>



**Key Idea: Reorder data to reduce the number of transmitted 1s**



**Model data reordering as an instance of the Traveling Salesman Problem**



- Unencoded
- Base+XOR
- Base+XOR with CDR (Best Stride)

**22.4% reduction in DRAM IO energy over Base+XOR Encoding**

UNIVERSITY of ROCHESTER