# Rethinking Graph Data Placement for Graph Neural Network Training on Multiple GPUs
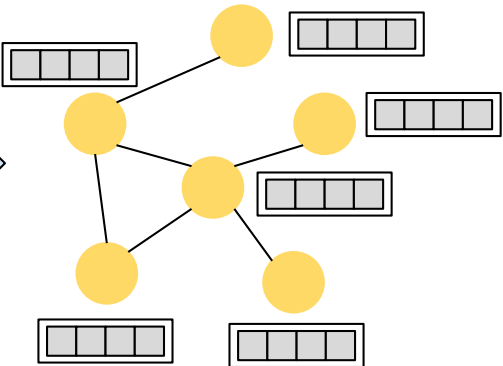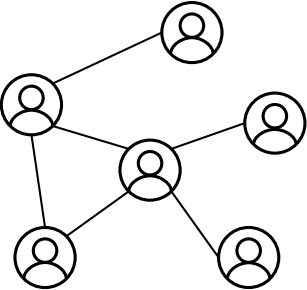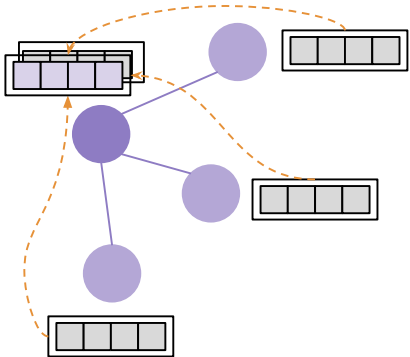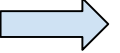
**Shihui Song and Peng Jiang**
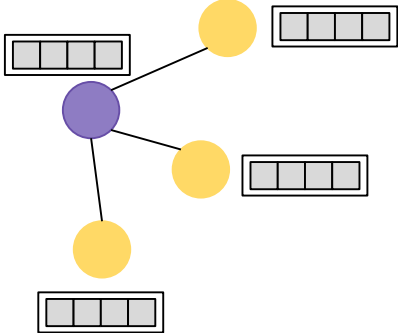
University of Iowa

# Graph Neural Network (GNN)



Graph Construction

Feature Aggregation

Graph Neural Network Model

Input

Hidden layer

Activation Function

Hidden layer

Activation Function

Output

# Sampling-based GNN training



- To reduce the computation, sampling-based GNN training samples a subset of neighbors and estimates the aggregation results based on the sampled nodes.

# Analysis of DGL and PaGraph (PG)



**DGL**
- Adopts METIS graph partitioning
- Assumes that the graph can be entirely stored on multiple GPUs

**PaGraph**
- Stores the graph on CPU and buffers the most frequently accessed nodes of each partition on GPU

# Motivation

- Assume that the GPU memory is small and we can only store 20% of nodes that are most frequently accessed on each GPU



Loading features is a bottleneck of training

# Performance Model



Goal ➤ Find the optimal configuration of $B_i$

**Cost Function of GPU_i:**

$$C_i(B) = \begin{cases} C_{cpu}\left|R_i \setminus B_{gpu}\right| + C_{gpu}\left|R_i \cap B_{gpu}\right|, & \text{if } C_{gpu} < C_{cpu}; \\ C_{cpu}\left|R_i\right| & \text{if } C_{gpu} \geq C_{cpu}. \end{cases}$$

$R_i$ : The nodes GPU_i needs to read from CPU and other GPUs

$B_{gpu}$: The nodes saved on GPUs

$C_{cpu}$ : The cost of reading a node on CPU

$C_{gpu}$ : The cost of reading a node on a different GPU

**Optimization Problem:**

$$\min \quad \max_{i \in [1,n]}(\mathrm{E}_{S_i \sim D}[C_i(B)]),$$

$$\text{subject to} \quad \left|B_i\right| \leq BSIZE, \quad i = 1, \dots, n$$

# Case1 $C_{gpu} \geq C_{cpu}$

**Optimization Problem:**

$$\mathrm{E}_{S_i}[C_i] = C_{cpu}\mathrm{E}_{S_i}[|R_i|]$$

$$= C_{cpu}\mathrm{E}_{S_i}[|S_i \setminus B_i|]$$

$$= C_{cpu}\mathrm{E}_{S_i}[|S_i|] - C_{cpu}\mathrm{E}_{S_i}[|S_i \cap B_i|]$$

$S_i$ : The nodes GPU_i needs to read

$B_i$ : The nodes saved on GPU_i local memory

When $\mathrm{E}_{S_i}[|S_i \cap B_i|]$ is maximized for every GPU_i, this formular can achieve the minimum value.

Rule1: We store nodes with the highest sampling probability on it

# Case2 $C_{gpu} < C_{cpu}$

**Optimization Problem:**

$$\mathrm{E}_{S_i}[C_i]$$

$$= C_{cpu}\mathrm{E}_{S_i}\left[\left|R_i \setminus B_{gpu}\right|\right] + C_{gpu}\mathrm{E}_{S_i}\left[\left|R_i \cap B_{gpu}\right|\right]$$

$$= C_{cpu}\mathrm{E}_{S_i}\left[\left|R_i\right|\right] - \left(C_{cpu} - C_{gpu}\right)\mathrm{E}_{S_i}\left[\left|R_i \cap B_{gpu}\right|\right]$$

$$= C_{cpu}\mathrm{E}_{S_i}\left[\left|S_i\right|\right] - \left(C_{cpu} - C_{gpu}\right)\mathrm{E}_{S_i}\left[\left|S_i \cap B_{gpu}\right|\right] - C_{gpu}\mathrm{E}_{S_i}\left[\left|S_i \cap B_i\right|\right]$$

Tradeoff

Should store as many nodes as possible on all GPUs

Each GPU stores the same set of nodes with the highest sampling probability

# Case2 Algorithm

**Algorithm 1:** Distributing node features onto multiple GPUs with fast interconnects

**Input:** $\alpha$; #nodes: $N$; #GPUs: $n$; Sampling probability: $p$; Buffer size: BSIZE

**Output:** $B = \{B_1, \ldots, B_n\}$

/* Sort nodes by probability $p$ in descending order    */

1  $V = \text{sort\_nodes}(N, p)$;

/* Initialize buffer on each GPU with nodes of highest sampling probabilities    */

2  **for** $i = 1$ **to** $n$ **do**

3      $B_i = [V[0], V[1], \ldots V[BSIZE - 1]]$;

4  $p\_sum = [0.0, \ldots, 0.0]$;

Increase of the second term  $C_{gpu}\left(p(old\_node) - p(new\_node)\right)$

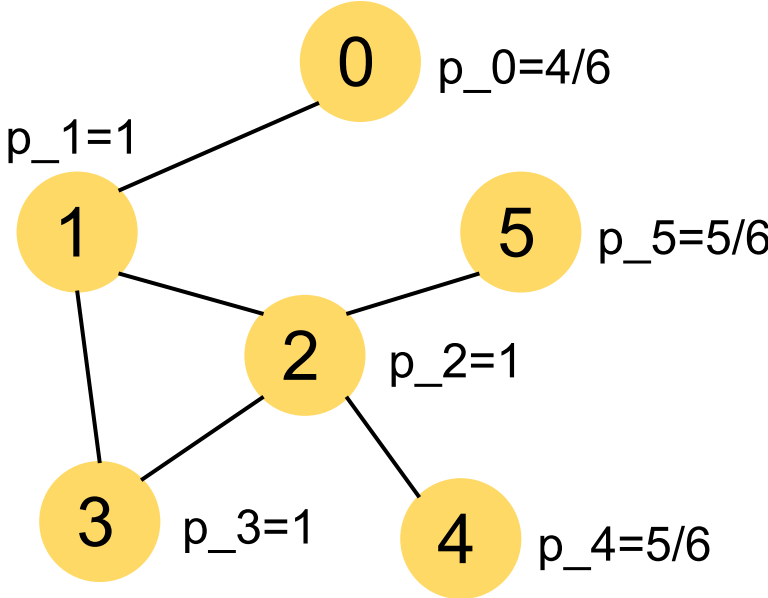Decrease of the last term  $\left(C_{cpu} - C_{gpu}\right) p(new\_node)$

Beneficial condition  $p(new\_node) > \dfrac{C_{gpu}}{C_{cpu}} p(old\_node)$

5  **for** $i = 0$ **to** $(\min(N, n \cdot BSIZE) - BSIZE - 1)$ **do**

6      **if** $i \bmod n == 0$ **then**

      /* Sort GPUs by $p\_sum$ in ascending order    */

7        $ordered\_devices = \text{sort\_device}(n, p\_sum)$;

    /* Do not change last device in each round    */

8      **if** $i \bmod n == n - 1$ **then continue**;

    /* Get device in the sorted order    */

9      $device = ordered\_devices[i \bmod n]$;

    /* Select the next node in $V$    */

10      $new\_node = V[i + BSIZE]$ ;

    /* Select a duplicate node on device    */

11      $old\_node\_idx = BSIZE - 1 - \lfloor i/n \rfloor$;

12      $old\_node = V[old\_node\_idx]$;

    /* Check if the replacement is beneficial    */

13      **if** $p_{(new\_node)} > \alpha \cdot p_{(old\_node)}$ **then**

      /* Replace $old\_node$ on device with $new\_node$    */

14        $B_{device}[old\_node\_idx] = new\_node$;

      /* Update $p\_sum$    */

15        $p\_sum[device] += p_{(new\_node)}$

16      **else break**;

17  **return** $\{B_1, \ldots, B_n\}$

# Case2 Example



The ordered nodes: [1, 2, 3, 4, 5, 0]
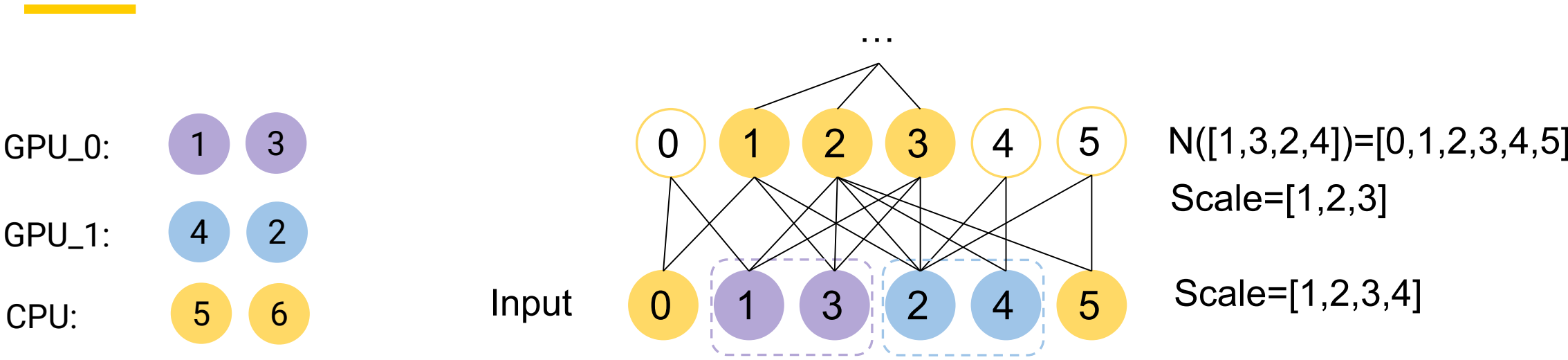
# Limitation of data placement

If the access frequency is less skewed and the GPU memory is small, the data loading might be expensive even with algorithm1

- We still need to load feature embeddings from CPU in most of the cases

→ A locality-aware neighbor sampling technique to further reduce the data movement overhead

# Locality-Aware Neighbor Sampling

GPU_0: 1 3

GPU_1: 4 2

CPU: 5 6

...

0 1 2 3 4 5

N([1,3,2,4])=[0,1,2,3,4,5]

Scale=[1,2,3]

Input 0 1 3 2 4 5

Scale=[1,2,3,4]

The ordered nodes: [1, 2, 3, 4, 5, 0]

Multiply the sampling probabilities of the neighbor set of B with an adjustable factor

1 2 4 … Maximum

# Experimental Setup

- Platform

  A single machine with two Intel Xeon Gold 6248 CPUs and eight Nvidia Tesla V100 GPUs
  - GPUs connected with NVLink Bridge: (2+2)GPU, (2*4)GPU
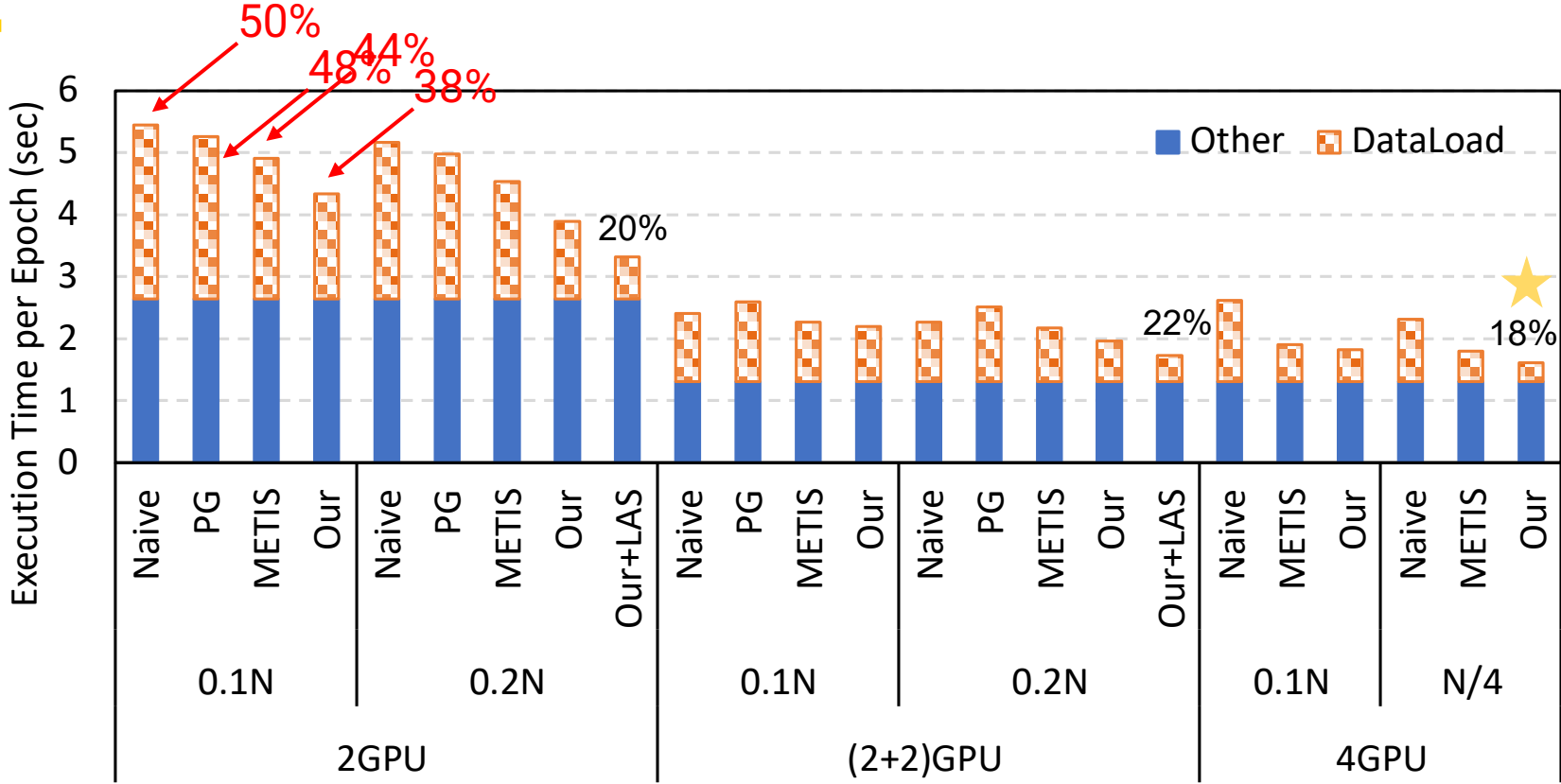  - GPUs connected with NVSwitch: 2GPU, 4GPU, 8GPU

- Dataset

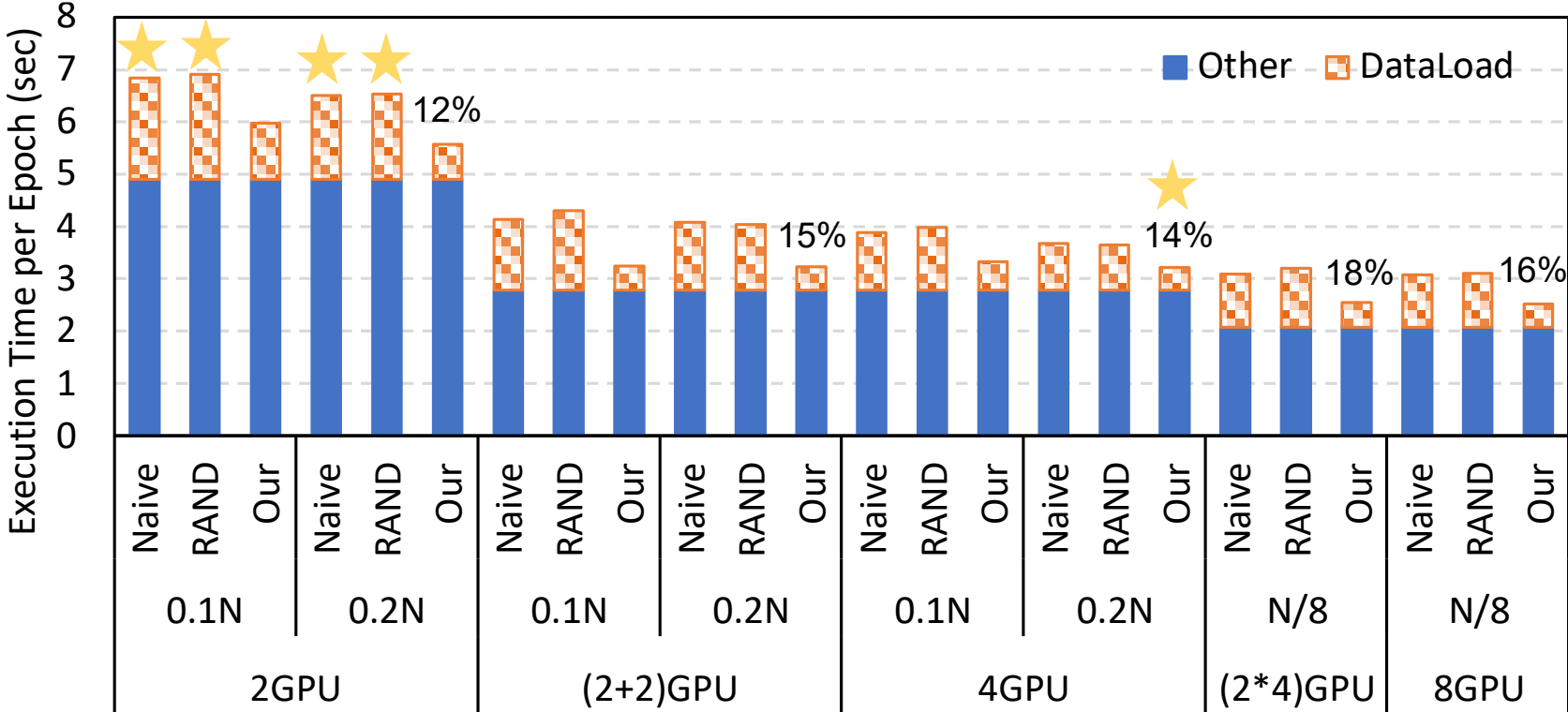| | Reddit | Yelp | Products | Papers100M | MAG240M |
|---|---|---|---|---|---|
| #nodes | 233K | 717K | 2.4M | 111M | 122M |
| #edges | 11.6M | 7.0M | 62M | 1.6B | 1.3B |
| feat_size | 535MB | 820MB | 934MB | 53GB | 175GB |

- Baseline
  - Naive partitioning
  - Random partitioning
  - METIS partitioning
  - PaGraph partitioning (Lin et al., SoCC'20)

# Evaluation: Speedup on Reddit



Execution Time per Epoch (sec)

Legend: Other, DataLoad

50% 48% 44% 38% 20% 22% 18%

Categories (left to right):
- 2GPU, 0.1N: Naive, PG, METIS, Our
- 2GPU, 0.2N: Naive, PG, METIS, Our, Our+LAS
- (2+2)GPU, 0.1N: Naive, PG, METIS, Our
- (2+2)GPU, 0.2N: Naive, PG, METIS, Our, Our+LAS
- 4GPU, 0.1N: Naive, METIS, Our
- 4GPU, N/4: Naive, METIS, Our

**Experiments**

# Evaluation: Speedup on Papers100M

# Evaluation: Accuracy



Training loss on Reddit (2GPU)

Training loss on Reddit (2+2GPU)

LAS has similar convergency speed

LAS has smaller loss at the end of training

# Evaluation: Preprocess overhead

Our algorithm is much faster than the previous graph partitioning algorithms

The execution time for dividing the graphs into four parts

|          | Reddit | Yelp | Products |
|----------|--------|------|----------|
| PaGraph  | 382    | 1976 | 4753     |
| METIS    | 17     | 15   | 83       |
| Our      | 0.49   | 0.76 | 3.6      |

- PaGraph has $O(N^2)$ time complexity
- We have $O(N)$ time and space complexity

# Summary

- Aim to reduce the data loading overhead for largescale GNN training on multiple GPUs

- Propose a performance model of the data movement among CPU and GPUs and provide an efficient algorithm to find an optimal data placement strategy

- Propose a locality-aware neighbor sampling technique to further reduce the data loading overhead

- Reduce data movement overhead by 1.2x to 3.3x times with data placement strategy, and achieve up to 4.4x times speedup with locality-aware sampling

For information, doubts and clarifications, contact: shihui-song@uiowa.edu