# LOW-RANK OPTIMIZATION FOR DATA SHUFFLING IN WIRELESS DISTRIBUTED COMPUTING

Kai Yang\*, Yuanming Shi\*, and Zhi Ding†

\*School of Information Science and Technology, ShanghaiTech University, Shanghai, China †Dept. of ECE, University of California, Davis, California 95616, USA

#### On-Device Intelligence

Deploy artificial intelligence on mobile devices: augmented reality, smart vehicles, drones etc.

Benifits: smart decision, low latency, privacy

**Key challenges:** performing intensive computation tasks on resource-constrained mobile devices

**Solution:** Wireless distributed computing – pooling computation & storage resources among the devices

## Wireless Distributed Computing: Computation Framework

**Target:** K users with individual computation task  $\phi_k$ , depending on a common dataset

Dataset: too large for storing in one mobile device (e.g., a feature library of objects)

**Solution:** stored across devices (separated to N files  $\{f_1,\cdots,f_N\}$ , but each can only store up to  $\mu$ files), supported by popular distributed computing framework **MapReduce** 

$$\phi_k(d_k;f_1,\cdots,f_N)=h_k(g_{k,1}(d_k;f_1),\cdots,g_{k,N}(d_k;f_N)).$$
 Map function:  $w_{k,t}=g_{k,t}(d_k;f_t)$ ; Reduce function:  $h_k(w_{k,1},\cdots,w_{k,N})$ ; Input:  $d_k$ 

# Overall procedure for distributed

## computing:

• Dataset Placement Phase: determine file placement strategy index set of files stored at each node  $\mathcal{F}_k \subseteq [N]$ ) and delivery files in advance

Map Phase: compute intermediate values  $\{w_{s,t}:s\in[K],t\in\mathcal{F}_k\}$  locally for all  $k\in[K]$ 

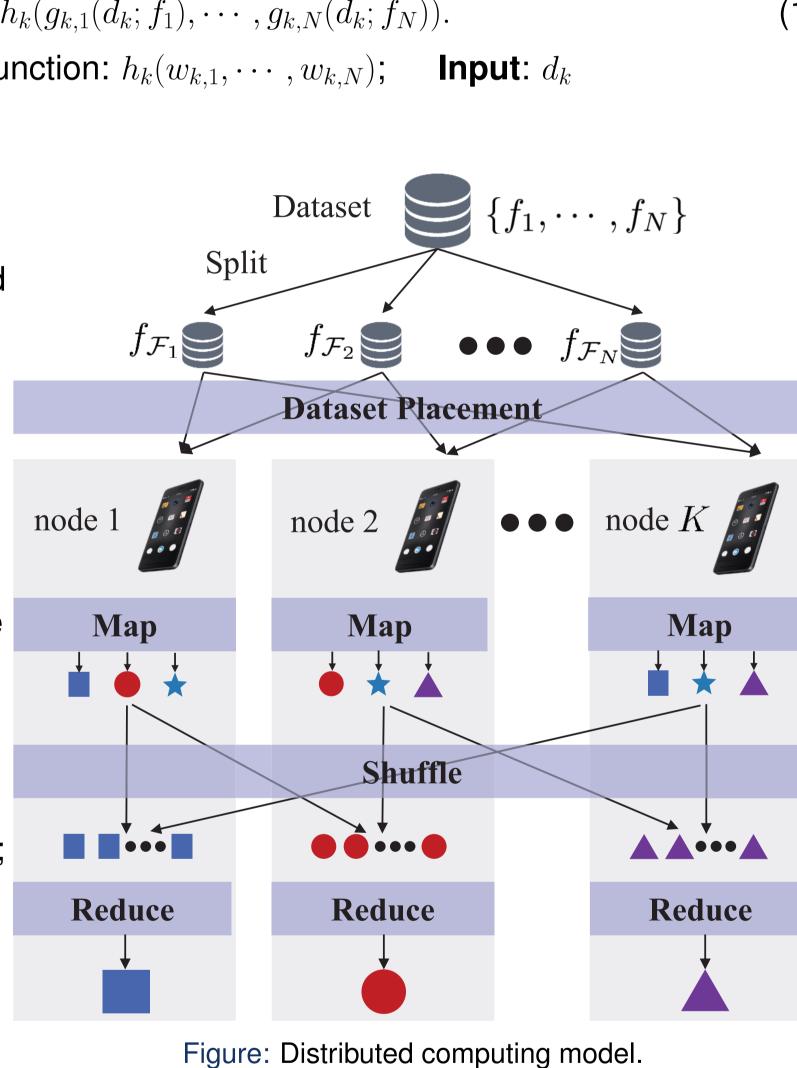
 Shuffle Phase: exchange intermediate values wirelessly among nodes.

 Reduce Phase: construct the output value using the reduce function

# Assumptions:

 $h_k(w_{k,1},\cdots,w_{k,N})$  .

 $\mu < N$ ; the sizes of inputs and intermediate values are small enough to be stored locally; overhead of collecting all inputs is negligible



#### Challenge and Prior Works

Technical challenge: To enable real-time and low-latency applications, inter-device communications for data shuffling and scheduling become the main bottleneck

Prior works: [Li, et al., TIT 18] [Li, et al., TON 17] proposed coding schemes to reduce communication load (e.g., information bits) for data shuffling in wireline and wireless distributed computing systems; orthogonal uplink transmission in [Li, et al., TON 17]

But the communication efficiency (i.e., achieved data rates) for data shuffling is also critical in wireless networks due to the limited spectral resources and interference!

Our Work: Low Rank Optimization for Data Shuffling

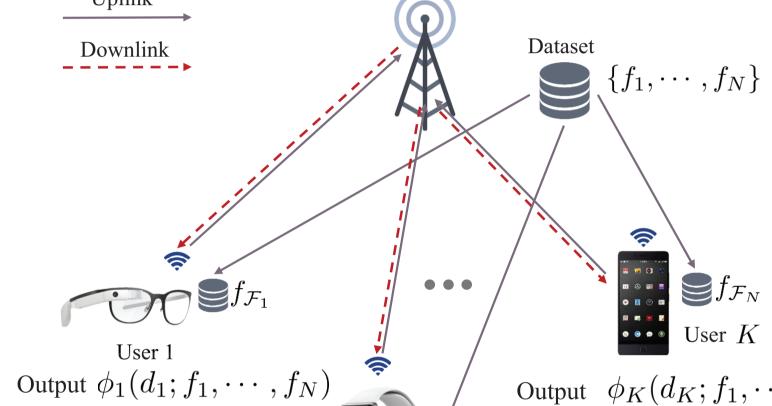
#### **Communication Model**

**Network model**: K users (each with L antennas) exchange intermediate values via a wireless access point (M antennas); **co-frequency transmission** in both uplink and downlink

Shuffle Phase reformulated as a message delivery problem:

- entire set of messages  $\{W_1, \cdots, W_T\}$  (intermediate values $\{w_{k,n}\}$ , T = KN)
- index set of messages available (can be computed locally) at user k:  $\mathcal{T}_k \subseteq [T]$
- index set of messages required by mobile user k:  $\mathcal{R}_k \subseteq [T]$

# Uplink multiple access (MAC) stage:



User 2 Figure: Communication model.

# $oldsymbol{y} = \sum (oldsymbol{H}_k^u \otimes oldsymbol{I}_r) oldsymbol{x}_k + oldsymbol{n}^u,$

 $\boldsymbol{x}_k \in \mathbb{C}^{Lr}$ : transmitted by mobile user kr: channel uses (block fading channel) **Downlink broadcasting (BC)** stage:

 $oldsymbol{y} \in \mathbb{C}^{Mr}$ : received at the AP

 $oldsymbol{z}_k = (oldsymbol{H}_k^d \otimes oldsymbol{I}_r) oldsymbol{y} + oldsymbol{n}_k^d.$  $\boldsymbol{z}_k \in \mathbb{C}^{Lr}$ : received by mobile user kOverall input-output relationship from mobile user to mobile user:

 $oldsymbol{z}_k = \sum (oldsymbol{H}_{ki} \otimes oldsymbol{I}_r) oldsymbol{x}_i + oldsymbol{n}_k,$  (4)  $oldsymbol{H}_{ki} = oldsymbol{H}_k^d oldsymbol{H}_i^u \in \mathbb{C}^{L imes L}$  ,

**Performance metric:** For message l at user k, the degree-of-freedom (**DoF**) defined by

$$\operatorname{DoF}_{k,l} \stackrel{\triangle}{=} \lim \sup_{\operatorname{SNR}_{k,l} \to \infty} \frac{R_{k,l}}{\log(\operatorname{SNR}_{k,l})}.$$

 $oldsymbol{n}_k = (oldsymbol{H}_k^d \otimes oldsymbol{I}_r) oldsymbol{n}^u + oldsymbol{n}_k^d$ 

# Linear Coding Scheme with Interference Alignment

Precoding matrix  $m{V}_{kj}\in\mathbb{C}^{Lr imes d}$ :  $m{x}_k=\sum_{j\in\mathcal{T}_k}m{V}_{kj}m{s}_j$  Decoding matrix  $m{U}_{kl}\in\mathbb{C}^{d imes Lr}$ 

$$egin{aligned} ilde{m{z}}_{kl} &= m{U}_{kl} m{z}_k = m{U}_{kl} \sum_{i=1}^K (m{H}_{ki} \otimes m{I}_r) \sum_{j \in \mathcal{T}_i} m{V}_{ij} m{s}_j + m{n}_{kl} \ &= \mathcal{I}_1(m{s}_l) + \mathcal{I}_2(m{s}_j : j \in \mathcal{T}_k) + \mathcal{I}_3(m{s}_j : j \notin \mathcal{T}_k \cup \{l\}\}) + ilde{m{n}}_{kl} \ & ext{desired message} & ext{locally available messages} & ext{interferences} \end{aligned}$$

 $egin{aligned} \mathcal{I}_1(oldsymbol{s}_l) &= \sum_{i:l \in \mathcal{T}_i} oldsymbol{U}_{kl}(oldsymbol{H}_{ki} \otimes oldsymbol{I}_r) oldsymbol{V}_{il} oldsymbol{s}_l \end{aligned} \qquad \mathcal{I}_2(\{oldsymbol{s}_j: j \in \mathcal{T}_k\}) = \sum_{j \in \mathcal{T}_k} \sum_{i:j \in \mathcal{T}_i} oldsymbol{U}_{kl}(oldsymbol{H}_{ki} \otimes oldsymbol{I}_r) oldsymbol{V}_{ij} oldsymbol{s}_j \end{aligned}$  $\mathcal{I}_3(\{oldsymbol{s}_j:j
otin\mathcal{T}_k\cup\{l\}\})=\sum_{j
otin\mathcal{T}_k\cup\{l\}}\sum_{i:j\in\mathcal{T}_i}oldsymbol{U}_{kl}(oldsymbol{H}_{ki}\otimesoldsymbol{I}_r)oldsymbol{V}_{ij}oldsymbol{s}_j.$ 

#### If interference alignment conditions:

$$\det \left( \sum_{i:l \in \mathcal{T}_i} \mathbf{U}_{kl} (\mathbf{H}_{ki} \otimes \mathbf{I}_r) \mathbf{V}_{il} \right) \neq 0,$$

$$\sum_{i:j \in \mathcal{T}_i} \mathbf{U}_{kl} (\mathbf{H}_{ki} \otimes \mathbf{I}_r) \mathbf{V}_{ij} = \mathbf{0}, \ j \notin \mathcal{T}_k \cup \{l\}$$

are satisfied, symmetric DoF  $\mathrm{DoF}_{\mathsf{sym}} = d/r$ is achievable

Assume  $\sum_{i:l\in\mathcal{T}_i}oldsymbol{U}_{kl}(oldsymbol{H}_{ki}\otimesoldsymbol{I}_r)oldsymbol{V}_{il}=oldsymbol{I}$  w.l.o.g.

Low rank formulation: D = LdKT

Nuclear norm relaxation yields poor performance due to the poor stucture of A. E.g., in a 2-user case K=N=2,  $\mu=d=L=M=1$ ,  $\mathscr P$  becomes

# Low-Rank Matrix Representation:

$$egin{aligned} oldsymbol{U}_{kl}(oldsymbol{H}_{ki}\otimes oldsymbol{I}_r)oldsymbol{V}_{ij} &= \sum_{m=1}^L \sum_{n=1}^L H_{ki}[m,n]oldsymbol{U}_{kl}[m]oldsymbol{V}_{ij}[n] \ oldsymbol{X} &= [oldsymbol{U}_{kl}[m]oldsymbol{V}_{ij}[n]] = egin{bmatrix} oldsymbol{U}_{11}[1] \ oldsymbol{U}_{KT}[L] \end{bmatrix} egin{bmatrix} oldsymbol{V}_{11}[1] & \cdots & oldsymbol{V}_{KT}[L] \end{bmatrix} \\ &\operatorname{rank}(oldsymbol{X}) &= r, oldsymbol{X} \in \mathbb{C}^{LdKT \times LdKT} \end{aligned}$$

 $\mathscr{P}: \underset{oldsymbol{X} \in \mathbb{C}^{D imes D}}{\operatorname{minimize}} \ \operatorname{rank}(oldsymbol{X})$ subject to  $\mathcal{A}(\boldsymbol{X}) = \boldsymbol{b}$ 

> $\underset{\boldsymbol{x}}{\text{minimize }} \text{ rank}(\boldsymbol{X})$ subject to  $\mathbf{X} = \begin{bmatrix} \star & \star & 1/H_{12} & 0 \\ 0 & 1/H_{21} & \star & \star \end{bmatrix}$ , (8)

#### DC Algorithmic Approach

**Ky Fan** k-norm: [Watson, 1993] The Ky Fan k-norm of X = sum of largest-k singular values,

$$\|\mathbf{X}\|_k = \sum_{i=1}^k \sigma_i(\mathbf{X}),$$
 (9)

For any matrix  $X \in \mathbb{C}^{m \times n}$ ,

$$rank(\mathbf{X}) = \min\{k : \|\mathbf{X}\|_* - \|\mathbf{X}\|_k = 0, k \le \min\{m, n\}\}.$$
 (1

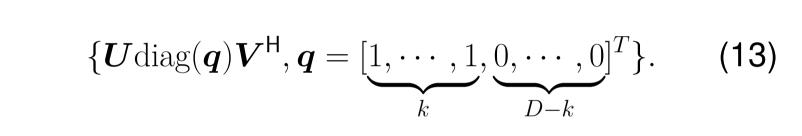
Then problem  $\mathscr{P}$  can be solved by finding the minimum k such that the optimal objective value is zero for the difference-of-convex-functions (DC)

minimize 
$$\|\boldsymbol{X}\|_* - \|\boldsymbol{X}\|_k$$
  
subject to  $\mathcal{A}(\boldsymbol{X}) = \boldsymbol{b}$ . (1)

Apply majorization-minimization (MM) algorithm to iteratively solve a convex approximation subproblem

minimize 
$$\|\boldsymbol{X}\|_* - \text{Tr}(\partial \|\boldsymbol{X}_{t-1}\|_k, \boldsymbol{X})$$
  
subject to  $\mathcal{A}(\boldsymbol{X}) = \boldsymbol{b}$ 

 $\partial |\!|\!| m{X}_t |\!|\!|_k$  is the subdifferential of  $|\!|\!| m{X} |\!|\!|_k$  at  $m{X}_t$  and given



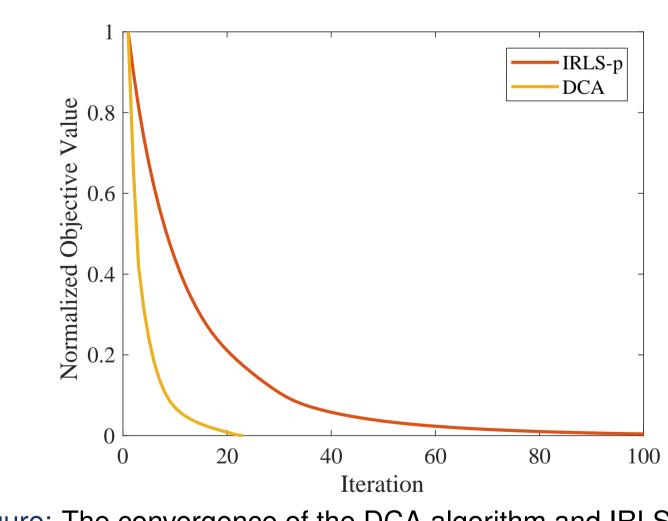
# Algorithm 1 DC Algorithm (DCA) for $\mathscr{P}$ Inputs: $\mathcal{A}, oldsymbol{b},$ , accuracy arepsilon2: Initialize: $\mathbf{s}_{:}$ for $k=1,\cdots,D$ do Initialize: $oldsymbol{X}_{0}^{[k]} \in \mathbb{C}^{D imes D}, t=1$ while not converge do Compute $\partial |\!|\!| oldsymbol{X}_{t-1}^{[k]} |\!|\!|_k$ Obtain the optimal solution $oldsymbol{X}_{\scriptscriptstyle t}^{[k]}$ of 12 end while if $\|m{X}^{[k]}\|_* - \|m{X}^{[k]}\|_k < arepsilon$ then return $m{X}^{[k]}$ 11: end for

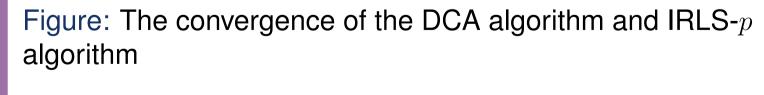
# Simulation

12: Outputs:

 $oldsymbol{X}^{[k]}$  and rank k .

Nuclear norm: nuclear norm relaxation approach. IRLS-p: iterative reweighted least square algorithm to further improve performance of nuclear norm relaxation. **DCA**: The present algorithm.





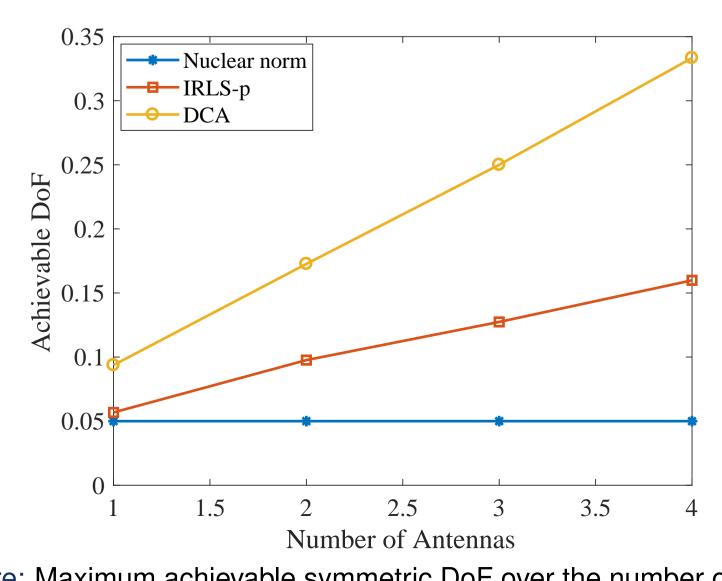
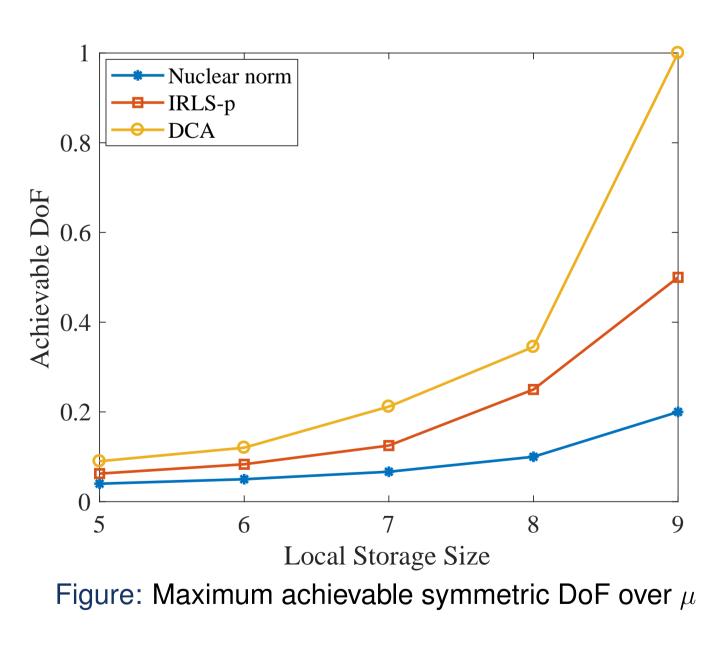


Figure: Maximum achievable symmetric DoF over the number of



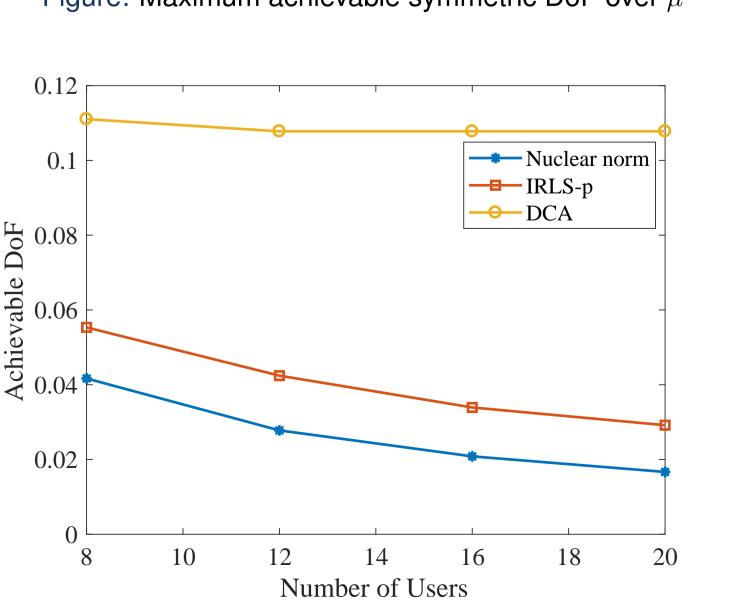


Figure: Maximum achievable symmetric DoF over the number of mobile

users with uniform placement strategy

## Conclusions

- Co-frequency transmission in both uplink and downlink
- Novel low-rank optimization approach for data shuffling in wireless distributed
- Algorithmic advantages of DC approach
- Scalability when the number of mobile users increases

# References

- [1] Song Han, Huizi Mao, and William J Dally, "Deep compression: Compressing deep neural networks with pruning, trained quantization and huffman coding," *Int. Conf. Learn. Representations (ICLR)*, 2016.
- [2] S. Li, Q. Yu, M. A. Maddah-Ali, and A. S. Avestimehr, "A scalable framework for wireless distributed computing IEEE/ACM Trans. Netw., vol. 25, no. 5, pp. 2643–2654, Oct. 2017.
- communication in distributed computing," *IEEE Trans. Inf. Theory*, vol. 64, no. 1, pp. 109–128, Jan. 2018. [4] GA Watson, "On matrix approximation problems with Ky Fan k norms," *Numerical Algorithms*, vol. 5, no. 5, pp.
- [5] Jun-ya Gotoh, Akiko Takeda, and Katsuya Tono, "DC formulations and algorithms for sparse optimization problems," Math. Program., to appear, 2017.

#### **Contact Information**

#### Authors:

- Kai Yang: yangkai@shanghaitech.edu.cn
- Yuanming Shi: shiym@shanghaitech.edu.cn
- Zhi Ding: zding@ucdavis.edu

ShanghaiTech University

UCDAVIS