



Computation Offloading in Mobile Edge Computing (MEC) Networks: A Multi-Task Learning Approach

Lijun Qian, Bo Yang, Joshua Bassey

CREDIT Center, Department of ECE, PVAMU

The 4th Mission Critical Big Data Analytics (MCBDA) Workshop

Contents

- 1 **Technical background**
- 2 **Proposed solution**
- 3 **Simulation results**
- 4 **Conclusions**



Cloud Computing (CC)



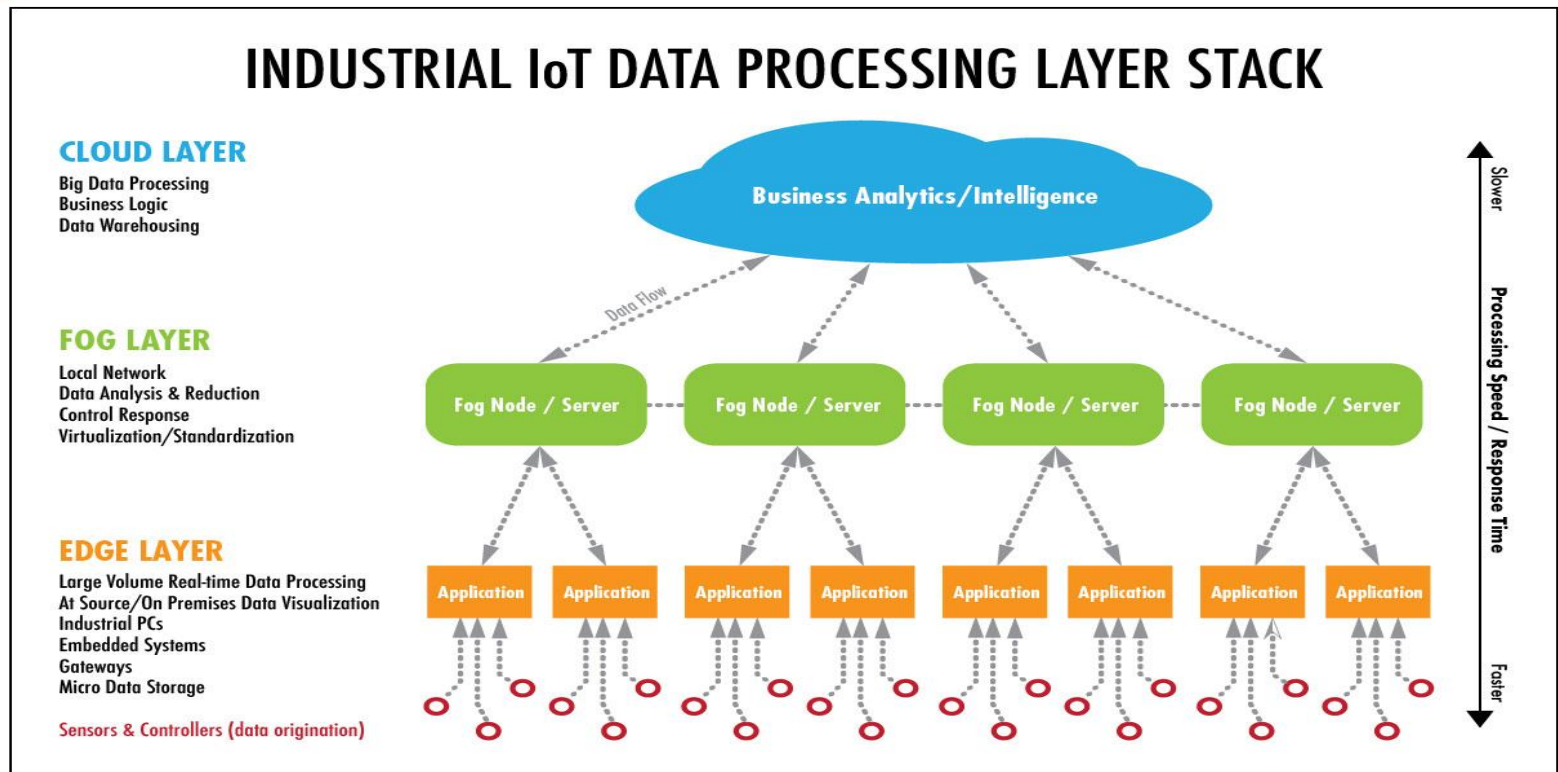
On-demand availability of computer system resources

The resources cannot be managed by the users

Wireless connection between users and central servers

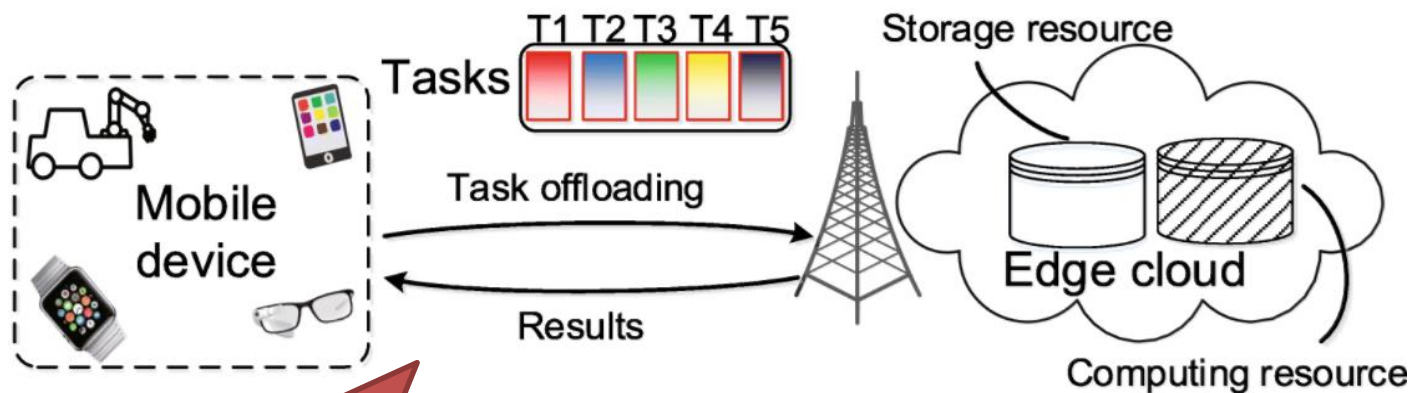
Mobile Edge Computing (MEC)

- MEC plays a key role in bringing cloud functionalities to the edge that **in close proximity to mobile users or devices**.



Computation Offloading

- Computation offloading is the transfer of resource intensive computational tasks to an external platform, such as a cluster, grid, or a cloud.
- Offloading may be necessary due to hardware limitations of a devices, such as limited computational power, storage, and energy.



Offloading decision includes:
1) Whether offloading or not?
2) How much resources can be obtained?

Offloading decision should be **OPTIMIZED** in **REAL-TIME!**

Challenges & motivations

- The authors generally formulated the joint optimization of computational resources and offloading decision as a **mixed integer nonlinear programming (MINLP) problem, which is NP-hard.**

Accuracy

- This NP-hard problem is usually solved by seeking an optimal (or sub-optimal) solution using traditional mathematical algorithms.

Efficiency

- Usually, huge computation overhead is introduced, which increases the delay obtaining the solutions.

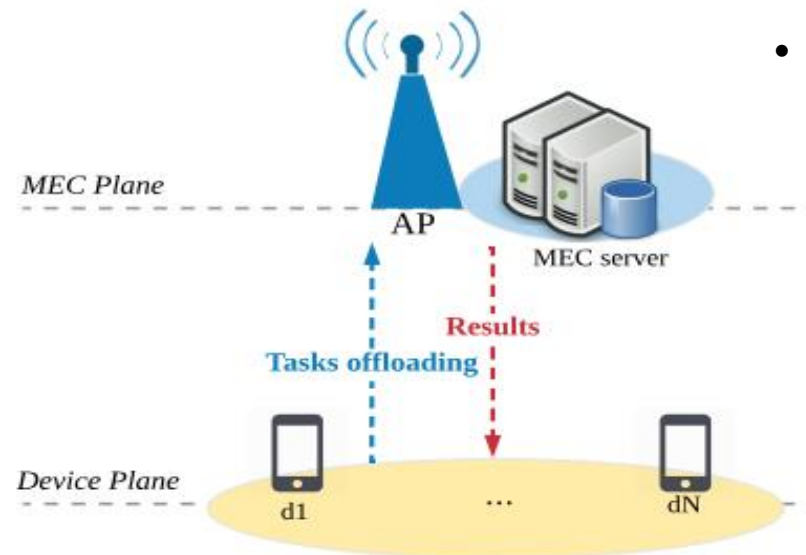
Contents

- 1 **Technical background**
- 2 **Proposed solution**
- 3 **Simulation results**
- 4 **Conclusions**



Network scenario & assumptions

- Total N devices in the network, $\mathcal{N} = \{d_1, d_2, \dots, d_N\}$.
- We assume that each device has only one computation-intensive task (denoted as $\mathcal{J}_i, \forall i \in \mathcal{N}$) to be processed at a time, which is atomic and cannot be further divided.
- Denote $D_i \in \{0, 1\}$ as the computation offloading decision of d_i , so the offloading decision vector for all devices can be expressed as $\mathbf{D} = \{D_i | i \in \mathcal{N}\}$, which is a N -dimensional binary vector.



- Let $\mathbf{F} = \{f_i | i \in \mathcal{N}\}$ be the allocated computational resource (i.e., central processing unit (CPU) cycles per second) vector by the MES. Hence, the offloading strategy for all devices can be defined as $\mathcal{S} = \{\mathbf{D}, \mathbf{F}\}$



System models

- Communication model:

- ❖ The mobile devices are assumed to operate using non-orthogonal multiple access (NOMA) such that the data from different mobile devices can be decoded separately from the superposed signal using multiuser detection algorithms at the AP.
- ❖ The received SINR of d_i , served by the AP is
$$\text{SINR}_i = \frac{P_t^i |h_i|^2}{\delta^2 + \sum_{j=i+1}^N P_t^j |h_j|^2}$$

- Computing model:

- ❖ Local execution delay: $\tau_l^i = \frac{c_i}{f_l^i}$
- ❖ Local energy consumption: $\varepsilon_l^i = \kappa (f_l^i)^2 c_i$
- ❖ Local weighted-cost: $\mathcal{O}_l^i = \alpha \tau_l^i + (1 - \alpha) \varepsilon_l^i, \forall i \in \mathcal{N}$
- ❖ Offloading delay and energy consumption:

$$\tau_o^i = \frac{s_i}{W \log_2(1 + \text{SINR}_i)} + \frac{c_i}{f_i}$$

$$\varepsilon_o^i = \frac{P_t^i s_i}{W \log_2(1 + \text{SINR}_i)} + \frac{P_f^i c_i}{f_i}$$

- ❖ Offloading weighted-cost:

$$\mathcal{O}_o^i = \alpha \tau_o^i + (1 - \alpha) \varepsilon_o^i, \forall i \in \mathcal{N}$$



Problem formulation and analysis

- We define the sum cost of the MEC system as the weighted-sum cost of all the devices: $O_{total} = \sum_{i \in \mathcal{N}} (1 - D_i) \mathcal{O}_l^i + D_i \mathcal{O}_o^i$
- We formulate the joint offloading and computational resource allocation as a weighted-sum cost minimization problem:

P1 (Original problem):

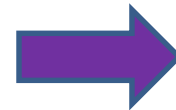
$$\underset{\{\mathbf{D}, \mathbf{F}\}}{\text{minimize}} \quad O_{total}$$

$$\text{s.t.} \quad \mathbf{C1} : D_i \in \{0, 1\}, \forall i \in \mathcal{N},$$

$$\mathbf{C2} : (1 - D_i) \tau_l^i + D_i \tau_o^i \leq \vartheta_i,$$

$$\mathbf{C3} : 0 \leq f_i \leq F, \forall i \in \mathcal{N},$$

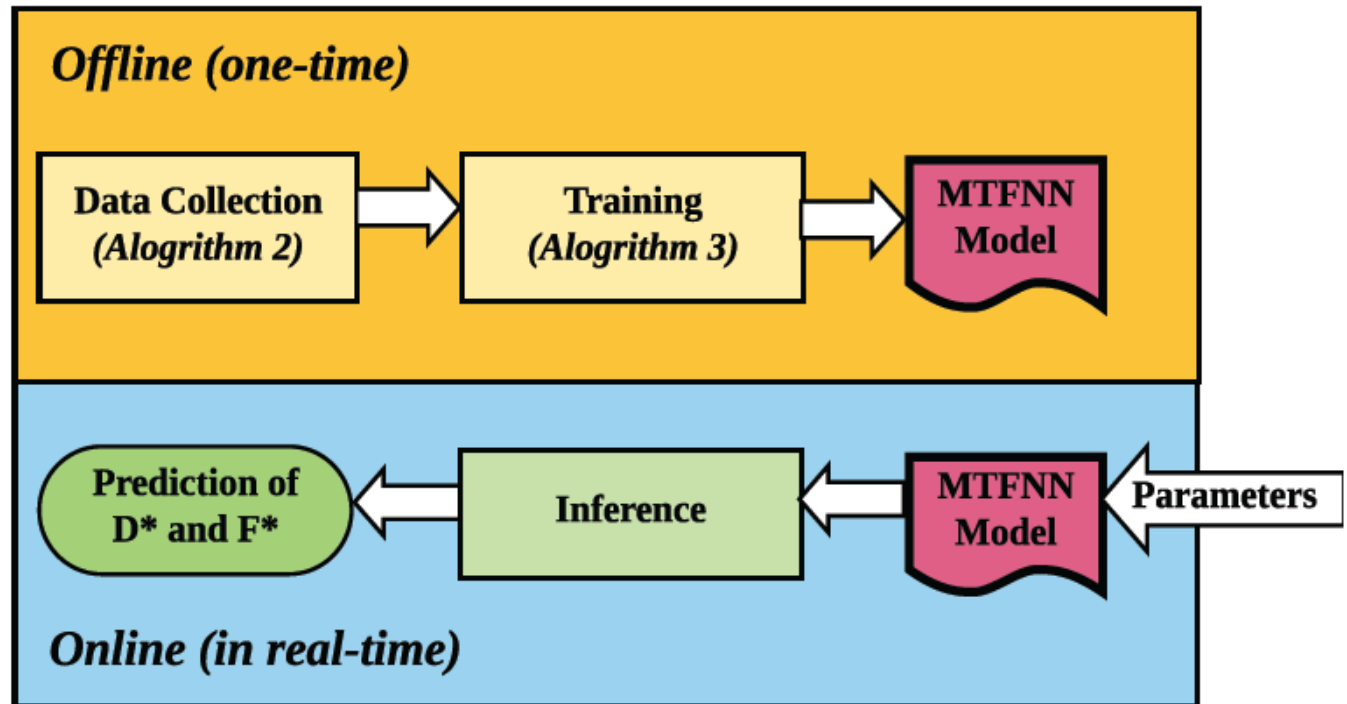
$$\mathbf{C4} : \sum_{i=1}^N D_i f_i \leq F, \forall i \in \mathcal{N}.$$



$$\mathbf{S}^* = \underset{\{\mathbf{D}, \mathbf{F}\}}{\text{argmin}} \quad O_{total}$$

Multi-task learning framework

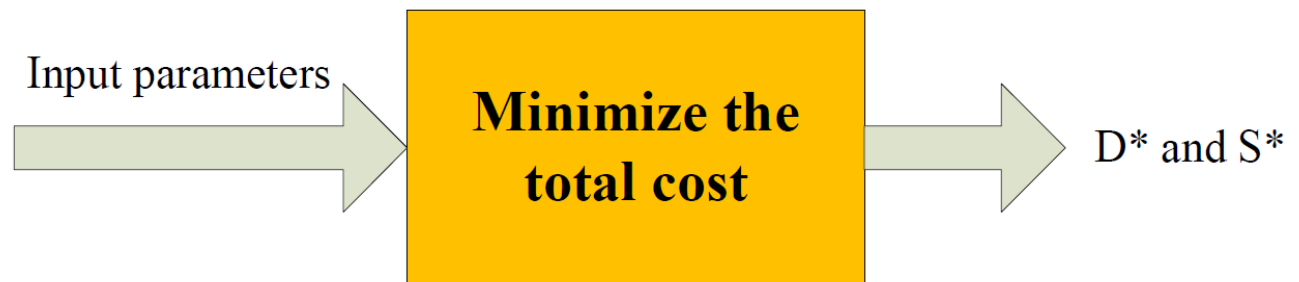
- The offloading decision making is formulated as a multiclass classification problem and the computational resource allocation is formulated as a regression problem.



Data collection (1/2)

- Given the parameters in the table, we generate the dataset by traversing all the possible combinations of **D** and **S** with the exhaustive searching algorithm, to minimize the weighted-sum cost.

Parameters	Value range
The number of devices (N)	[2 – 8]
Data payload size (s)	[1 – 500] kbits
CPU cycle required to process the data (c)	[3 – 1500] Megacycles
CPU frequency of the device (f_l)	[1Hz – 1GHz]
Weights of delay and energy cost (α, β)	[0.0 – 1.0]



Data collection (2/2)

- We generate and collect training dataset in the **MATLAB environment using a computer with NVIDIA GPU TITAN X (Pascal)**.
- The GPU can accelerate the matrix calculation in MATLAB.

Algorithm 2: Dataset Generation Algorithm

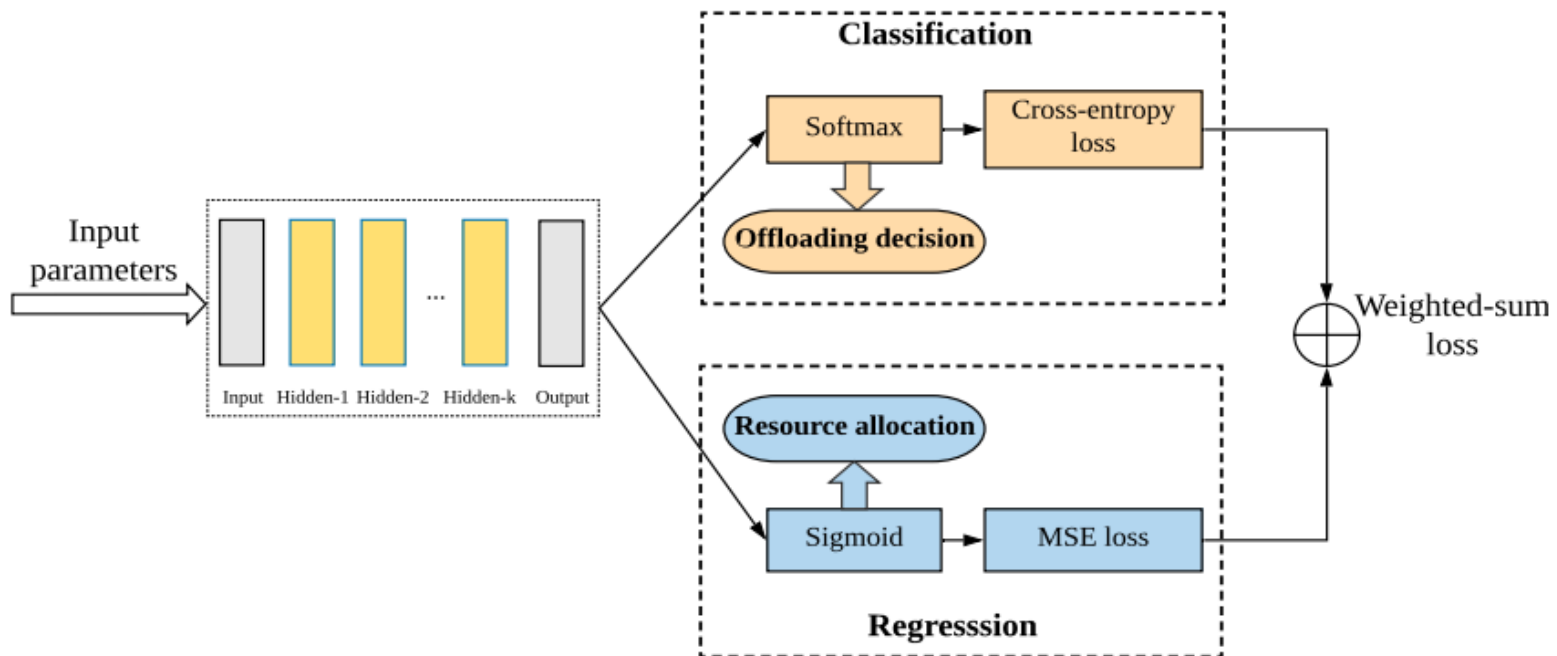
Initialization: $i = 0, S^* = \emptyset;$

Iteration:

- 1: **while** $i < \text{dataset size}$ **do**
 - 2: $i \leftarrow i + 1;$
 - 3: Generate input parameters set (X_i) for all devices;
 - 4: Formulate the optimization problem **P1** as (18);
 - 5: Solve **P1** with exhaustive searching method and record the optimal solution as $Y_i = (D^*, \Theta^*);$
 - 6: Add an input/output pair $S_i = \{X_i, Y_i\}$ to $S^*;$
 - 7: **end while**
-

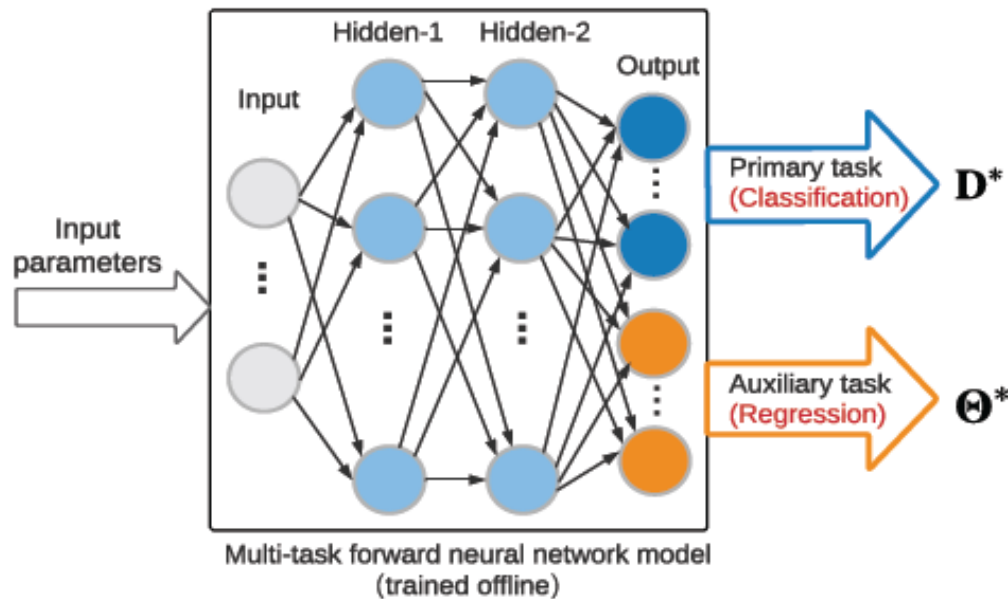
Offline training

- For the classifier, $l_c = -\frac{1}{K} \sum_{i=1}^K Y_i \ln f(\mathbf{X}_i)$,
 - For the regressor, $l_r = \frac{1}{n} \sum_{i=1}^n (Y_i - f(\mathbf{X}_i))^2$
- ➔ $l = \chi_c l_c + \chi_l l_r$



The pre-trained model

- An example of the pre-trained MTL model with $N = 3$



Contents

- 1 **Technical background**
- 2 **Proposed solution**
- 3 **Simulation results**
- 4 **Conclusions**



Testing results (1/2)

- We compare with a benchmark scheme “sBB” which is implemented using the MATLAB toolbox of the APMonitor Optimization Suite (<http://APMonitor.com>).
- Performance index:
 - Computation complexity:

$$t = \frac{\text{Total execution time}}{\text{Number of samples}}$$

The smaller
the better

- Computation accuracy:

Classification: $\eta = \frac{\text{Number of correct predictions}}{\text{Total number of predictions}}$

The larger
the better

Regression: $\varepsilon = \frac{1}{mN} \sum_{i=1}^m \sum_{j=1}^N (\bar{y}_j^i - x_j^i)^2$

The smaller
the better

Testing results (2/2)

η, ε, t \ S N	sBB	MTFNN with $\chi_c = \chi_l = 1$
2	70%, 0.055, 14.1 ms	96%, 0.016, 2.5 μ s
3	62%, 0.047, 14.2 ms	89%, 0.027, 2.5 μ s
4	58%, 0.053, 14.5 ms	83%, 0.029, 2.2 μ s
5	47%, 0.051, 15.2 ms	50%, 0.021, 2.0 μ s

η, ε, t \ S N	sBB	MTFNN with $\chi_c = 0, \chi_l = 1$
5	47%, 0.051, 15.2 ms	78%, 0.009, 5.0 μ s
6	42%, 0.092, 15.8 ms	82%, 0.009, 5.7 μ s
7	38%, 0.095, 16.6 ms	81%, 0.009, 3.6 μ s
8	34%, 0.097, 16.9 ms	78%, 0.009, 3.9 μ s



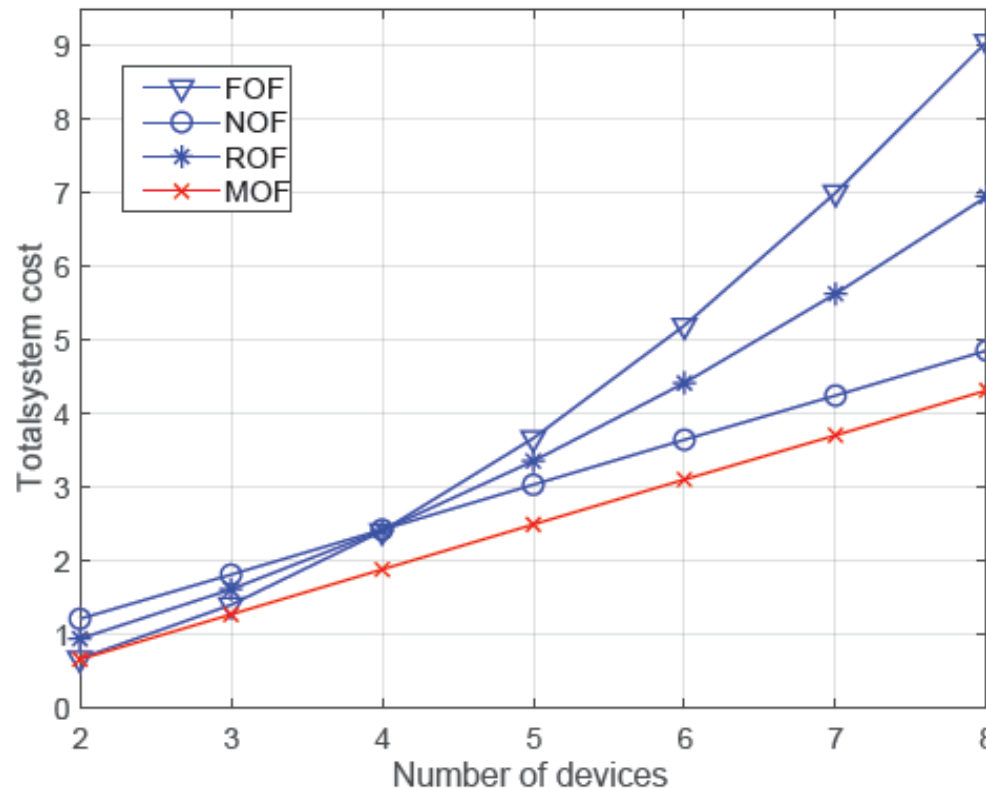
Simulation results (1/4)

- Three benchmark offloading approaches
 - Full offloading (FOF)
 - None offloading (NOF)
 - Random offloading (ROF): The ROF scheme denotes that all the tasks will be executed by the two ways above randomly.



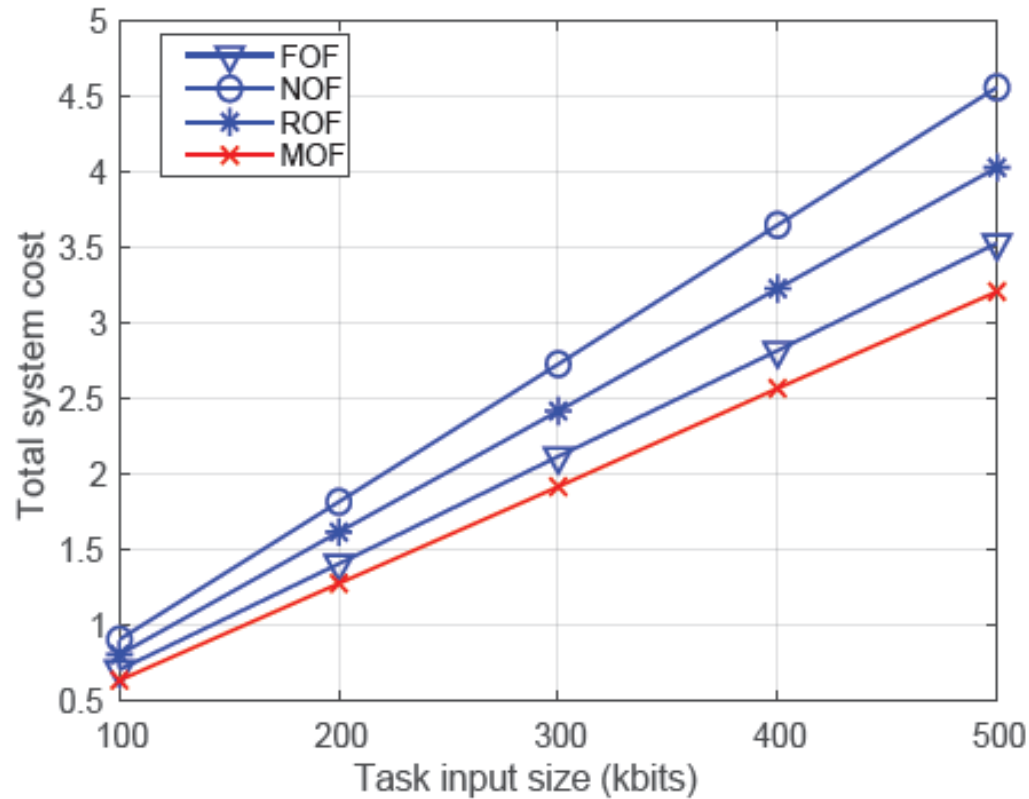
Simulation results (2/4)

$$f_l = 0.5 \text{ GHz} \quad \alpha = \beta = 0.5$$



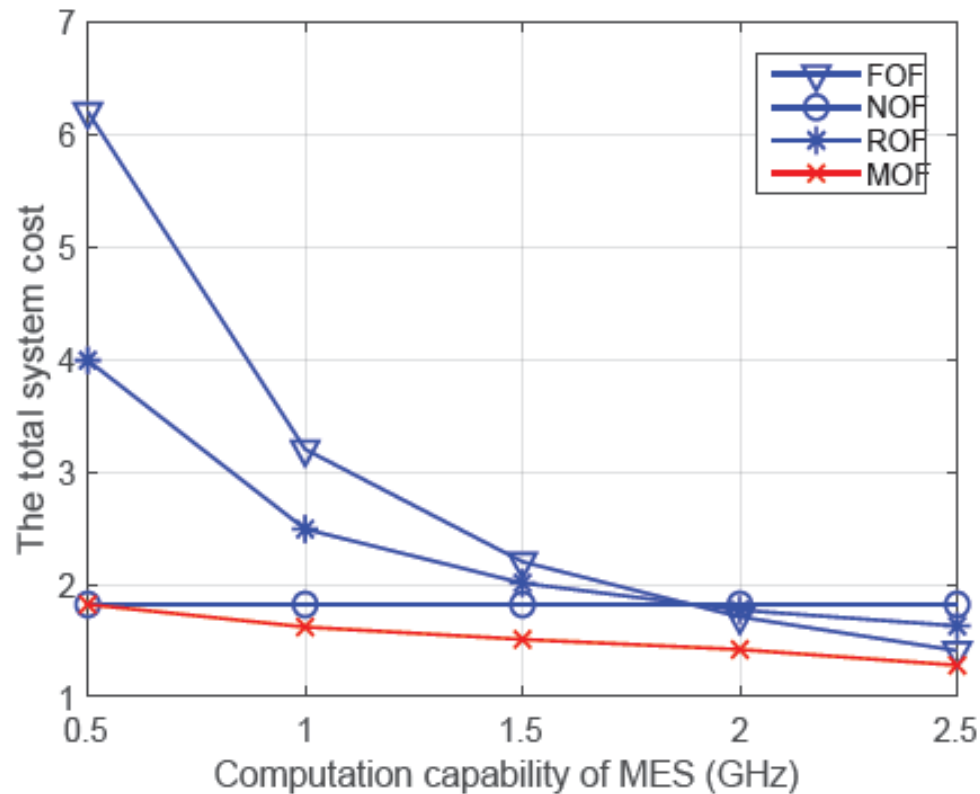
Simulation results (3/4)

$$f_l = 0.5 \text{ GHz} \quad \alpha = \beta = 0.5$$



Simulation results (4/4)

$$f_l = 0.5 \text{ GHz} \quad \alpha = \beta = 0.5$$



Contents

- 1 Technical background
- 2 Proposed solution
- 3 Simulation results
- 4 Conclusions



Contributions

- We propose a multi-task learning based solution that **can adapt to the varying network conditions and the changing requirements of devices' applications.**
- The MTL is **trained offline and only one time.** After the MTL model is trained, it can be directly used to generate the optimal solution of the MINLP problem with high accuracy in near-real-time.
- The proposed MTL model outperforms the conventional optimization algorithms significantly in terms of **computation time (four orders of magnitude)** and **inference accuracy (up to two times better).**



Q/A

