

Future Data-Aware Decision Making for Edge Computing

Dr Konstantinos (Kostas) Kolomvatsos

Monday, September 26th, 2022 School of Computing Science, University of Glasgow

SICSA DVF funding call







At a Glance

2016 Call Marie Skłodowska-Curie Action (MSCA) Individual Fellowship School of Computing Science University of Glasgow



2013 PhD in Computer Science National and Kapodistrian University of Athens

June 2020 Assistant Professor Department of Informatics and Telecommunications University of Thessaly http://kostasks.users.uth.gr

July 2020 Founder of the Intelligent Pervasive Systems (iPRISM) Research Group http://www.iprism.eu



Oct. 2020 Co-Founder of the Intelligent Systems for Orchestrating Pervasive Computing Applications (METIS) Research Lab http://metis.cs.uth.gr

Metis

Dec. 2020 Director of the METIS Lab



Current Activities:

- Applied Artificial Intelligence and Machine Learning
 - Distributed Intelligence
 - Pervasive Data Science



At a Glance

± iPRÍSM

Intelligent Pervasive Systems (iPRISM) http://www.iprism.eu

Lead: Dr Konstantinos (Kostas) Kolomvatsos

Research axes:

- Artificial Intelligence
- Applied (Deep) Machine Learning
- Computational Intelligence
- Distributed Intelligence
- Pervasive Computing
- Pervasive Data Science
- Proactive Decision Making
- Applications for Distributed Systems, Internet of Things, Edge Computing
- Predictive Intelligence
- Large Scale Data management





10/4/2022

Recent Research

01

02

03

04



University of Glasgow

Intelligent Systems in Pervasive, Edge Computing and Internet of Things

Contextual and Fuzzy Logic Reasoning for Pervasive Computing

Proactive Reasoning for Autonomous Behaviour and Decision Making

Pervasive Data Science Applications

Edge Computing

- Edge Computing (EC) deals with an additional infrastructure above the Internet of Things (IoT)
- EC 'imposes' an ecosystem of processing nodes that can execute tasks upon the collected data
- Gartner shared a report on ten (10) strategic trends affecting the Internet of Things (IoT) from 2019 to 2023 and beyond where the following are identified as the most impactful:
 - Artificial intelligence (AI)
 - The shift from intelligent edge to **intelligent mesh**
 - New IoT user experiences



Edge Computing

- We are at the early stages of the EC revolution to prepare the infrastructure for the new, modern, **Edge Mesh** (EM)
- EM provides a 'virtual' layer (<u>a computational/processing</u> <u>overlay</u>) that enables the cooperation between heterogeneous EC nodes to conclude a cooperative infrastructure close to end users
- Operators <u>can/should/will</u> open the ecosystem to third-parties, allowing them to rapidly deploy innovative applications and content



State of the Art

 \bigcirc

 Tasks are processed at the central server that will send the response back to the device

 \bigoplus

Cloud

 The whole process typically takes less than a second, however, there might be delays, e.g., due to a network glitch, weak internet connection, or long distance with the datacentre Every EC node is

[×≡ [

EC

responsible for processing and can respond to any request ✓ EC nodes are:
 resource constrained

EC Nodes

 \bigcirc

 \bigcirc

Ο

- o heterogeneous
- subjects to dynamic workloads change
- EM is a computation overlay network over the EC nodes

EM

- EM overcomes the problem of constrained resources through a cooperative model
- A mesh network of nodes enables distributed decisionmaking, sharing data and computation

Al/ML are required to support intelligence ✓ Distributed intelligence is the key enabler for cooperation and optimal decision making

 $\langle \rangle$

Intelligent EM

University of Glasgow

 \bigcirc

 \checkmark



How to define the network and computing model? How to distribute data processing?

How to jointly optimize computation? How to be stateful, i.e., exhibit different behaviour even for the same data according to the conditions met at a time instance?

University of Glasgow

he Scottish Informatics

Challenges





Modelling Challenges

Coordination Challenges

Support the creation/migration/replication of virtual resources at different levels of granularity

The demand for edge resources (even they are from peer nodes) have to be modelled (based on the traffic generated, this challenge is complex)

Mobility could result in the migration of services/data from one node onto another (Follow me Edge) or their replication (we define the Proactive Edge) Enable the required coordination that facilitates the management of geographically distributed devices

Facilitate fast migration/replication of abstract entities (such as functions or programs) or data

Evolution of a new model of multi-locational hybrid (edge & Cloud) data architectures

Exhibit the necessary intelligence for the proactive response to potential problems

Recent Publications



Kolomvatsos, K., 'A Proactive Inference Scheme for Data-Aware Decision Making in Support of Pervasive Applications', Future Generation Computer Systems (FGCS), Elsevier, 136, 2022, pp. 193-204

University of Glasgow



Kolomvatsos, K., Anagnostopoulos, C., 'A Proactive Statistical Model Supporting Services and Tasks Management in Pervasive Applications', IEEE Transactions on Network and Service Management, 2022, doi: 10.1109/TNSM.2022.3161663



Kolomvatsos, K., 'Data-driven Type-2 Fuzzy Sets for Tasks Management at the Edge', IEEE Transactions on Emerging Topics in Computational Intelligence, vol. 6, no. 2, pp. 377-386, April 2022

Research Axes









Research Axis A DATA





Data-aware Matching Inference

EC nodes, at regular intervals, exchange the calculated synopses

Without loss of generality, we consider that at t, a node n, receives N – 1 synopses {s^t_i} ∀j,j≠i.

n_i continuously monitors the discrepancy quanta with peers

The discrepancy quantum d^t_{ii} is calculated as the absolute value of the difference between the jth synopsis s^t and the local synopsis s^t

We generate the time series $d_{ij}^1, d_{ij}^2, \ldots, d_{ij}^W$ (sliding window) upon which the proposed 'inference process' is applied











Random Variables

S depicts the realization of the synopsis of a specific data dimension k

$$D_{ij} = \sum_{k=1}^{N} Z_k$$
$$Z_k = |S_{ik} - S_{ik}|$$

М

Expected Synopses Difference

Lemma. The expected difference between the synopses calculated by the ith node and the jth peer is given by $\mathbb{E}(D_{ij}) = \sum_{k=1}^{M} 2 \cdot A_k - \mu_{ik} - \mu_{jk}$ with $\mu_{ik} \otimes \mu_{jk}$ being the mean of the kth dimension in the ith and the jth synopses and $A_k = \int_{-\infty}^{+\infty} s \cdot [f_{S_{ik}}(s)F_{S_{ik}}(s) + f_{S_{ik}}(s)F_{S_{ik}}(s)] ds.$

Expected Difference of the Discrepancy Quanta

Proposition. The expected discrepancy quantum when data follow an Exponential distribution with the same rate λ is given by $\mathbb{E}(D_{ij}) = \frac{M}{\lambda}$.

Proposition. The expected discrepancy quantum when data follow an Exponential distribution with different rates is given by $\mathbb{E}(D_{ij}) = \sum_{k=1}^{M} \left(\frac{\lambda_{ik} + \lambda_{jk}}{\lambda_{ik} \lambda_{jk}} - \frac{2}{\lambda_{ik} + \lambda_{jk}} \right).$





University of Glasgow Data-aware Matching Inference

- ✓ The expected discrepancy quantum depicts the anticipated value of the difference between synopses at some point in the future
- \checkmark We combine such knowledge with the historical correlation of synopses to depict the trend of the discrepancy quanta
- \checkmark We adopt the known Pearson Correlation Coefficient (PCC) r_{μ}

$$r_{k} = \frac{\sum_{t=1}^{W} \left(s_{ik}^{t} - \mu_{ik} \right) \left(s_{jk}^{t} - \mu_{jk} \right)}{\sqrt{\sum_{t=1}^{W} \left(s_{ik}^{t} - \mu_{ik} \right)^{2}} \sqrt{\sum_{t=1}^{W} \left(s_{jk}^{t} - \mu_{jk} \right)^{2}}}$$

- ✓ Ideal scenario: Observe a positive correlation for the M dimensions
- ✓ **Real cases**: 'Mix' of positive or negative correlations
- \checkmark We assume a threshold θ_k over which we consider that the observed correlation is 'acceptable'
- ✓ We record the cardinality of the set $|\{r_k \ge \theta_k\}\forall k|$
- \checkmark The set of M indicators that should be aggregated into a single value
- ✓ We rely on a sparsity metric and define the logarithmic sparsity indicator $\rho \in R^+$ which depicts the population of unity values in

{r_k}, ∀k

$$\rho = \sum_{k=1}^{M} \log\left(1 + \mathbb{1}_{r_k \ge \theta_k}\right)$$

10/4/2022







Data-aware Matching Inference

- ✓ We create a temporal matching map and the selection of the appropriate peers for collaborative activities
- ✓ We define the *Matching Synopses Indicator* (MSI) **R** which aggregates
 - ✓ The expected discrepancy quantum $E(D_{ij})$
 - ✓ The correlation indicator ρ_{ij}
 - ✓ The communication cost c_{ij}

$$R_{ij} = \frac{e^{-\alpha \mathbb{E}(D_{ij})}}{1 + e^{-\beta \rho_{ij} + \gamma}} \frac{1}{1 + e^{-\delta c_{ij} + \epsilon}}$$

 $\alpha, \beta, \gamma, \delta, \varepsilon \in \mathbb{R}^+$ are smoothing parameters

- ✓ A sorted list $\{R_{ij}\}$ ∀j is provided in a descending order
- ✓ When required (e.g., to offload a task or 'borrow'/'lend' data from/to peers), every node can interact with peers exhibiting the highest R (a sub-set can be adopted)





Research Axis B SERVICES



Services Management Scenario



he Scottish Informatics

Services Management Scenario

We propose a model that deals with the decision of where to migrate a service The optimal migration strategy is intractable due to the dynamics of the EC ecosystem Tasks offloading can be affected by an additional level of decision and delays in the response Services migration should be carefully decided due to the resource constraints









Statistical Inference Utility based Decision Making Order statistics for analyzing Utility of the local presence of the demand of a service a service is compared to the utility of offloading tasks Target Aggregate a statistical inference technique with utility based decisions

10/4/2022

- ✓ { s_1 , s_2 , ..., s_N }: Set of services
- ✓ Services Demand Vector (SDV) for the ith node, i.e., d_i = {d_{i1}, d_{i2}, ..., d_{iN}}
- ✓ After the reception of a task T_{it}, we detect the required services and update the demand
- Our approach: keep the execution of popular tasks locally if the current load is at 'acceptable' levels (the node is not overloaded)

Decisions

Decision 1. Keep locally the execution of the task and, if needed, request the necessary services;Decision 2. Offload the task to the appropriate peer(s).



✓ d_i & l_i (load) are updated after the arrival of tasks
 ✓ We define the following variables:





✓ For both decisions, we define the expected utilities U & \hat{U} upon g and \hat{g}



10/4/2022

T Y O F AY BUS SA

- ✓ We assume that D_(r) is defined upon the random variable D with a random sample of size N and realizations **d** = [d₁, d₂, ..., d_N]
- ✓ For instance, $D_{(1)} = min(d_1, d_2, ..., d_N)$, $D_{(2)} = 2nd min(d_1, d_2, ..., d_N)$ and so on and so forth

Proposition. The expected utility for the local execution of a task that requires a service with demand d is given by $\mathbb{E}(G) = \frac{\varepsilon}{1+e^{-\gamma d+\delta}} F_{D_{(N-k)}}(d)$ where $F_{D_{(N-k)}}(d)$ is the cumulative distribution function (cdf) of the variable $D_{(N-k)}$. **Proposition.** The expected utility for the offloading action of a task that requires a service with demand d is given by $\mathbb{E}(\hat{G}) = \frac{\hat{\varepsilon}}{1+e^{-\hat{\gamma}d+\hat{\delta}}} \left(1-F_L(\hat{\theta})\right) \left(1-F_{D_{(N-k)}}(d)\right).$



Estimation of Load

- L is the random variable with realizations depicted by I
- ✓ Assume that nodes monitor L over the discrete time while storing the W recent values
- ✓ We adopt the widely known nonparametric Kernel Density Estimation (KDE) method for estimating the cdf and pdf of L

$$\hat{f}_L(l;W) = \frac{1}{W \cdot h} \sum_{j=1}^W K\left(\frac{|l - l_{t-W+j}|}{h}\right)$$
$$\hat{F}_L(l;W) = \frac{1}{W} \sum_{j=1}^W \frac{1}{2} \left(1 + \operatorname{erf}\left(\frac{l - l_{t-W+j}}{\sqrt{2}}\right)\right)$$



Theorem. The joint density function of $D_{(1)}, D_{(2)}, \ldots, D_{(N)}$ is given by $f_{1,2,\ldots,N}(d_1, d_2, \ldots, d_N) = N! f(d_1) f(d_2) \ldots f(d_N) \mathbb{I}_{d_1 < d_2 < \ldots < d_N}$.

Theorem. The probability of success for $D_{(r)}$ when the cdf of D_i is $F_D()$ is given by $F_{D_{(r)}}(x) = \sum_{j=r}^{N} {N \choose j} (F_D(x))^j (1 - F_D(x))^{N-j}$.

Proposition. The cdf of the N - k order statistic upon services demand values is given by $F_{D_{(N-k)}}(d) =$ $\sum_{j=N-k}^{N} {N \choose j} (F_D(d))^j (1-F_D(d))^{N-j}.$

Order

We can get that the pdf of the $D_{(r)}$

$$f_{D_{(r)}}(x) = \frac{N!}{(r-1)!(N-r)!} \left(F_D(x)\right)^{r-1} \left(1 - F_D(x)\right)^{N-r} f_D(x), x \in \mathbb{R}$$

University of Glasgow

Scenario A. Uniform distribution

$$f_{D_{N-k+1}}(x) = \frac{N!}{(N-k)!(k-1)!} x^{N-k} \left(1-x\right)^{k-1}$$

Scenario B. Exponential distribution

$$f_{D_{N-k+1}}(x) = \lambda \frac{N!}{(N-k)!(k-1)!} \left(1 - e^{-\lambda x}\right)^{N-k} e^{-\lambda kx}$$

Keep locally top-k services (based on demand)

Statistics Two scenarios



Algorithm Local Decision Making for t = 1, 2, ... do $\langle t, T_t, \mathscr{C}_t \rangle = getTask(\mathscr{T});$ Update(**d**); Calculate (g, \hat{g}) ; getExpectedDemandRankings(d); getExpectedUtilities($\mathbb{E}(G), \mathbb{E}(\hat{G})$); $U = \mathbb{E}\left(G\right) \cdot e^{-\eta \frac{\hat{\alpha}}{\zeta}}$ Calculate (U, \hat{U}) ; Decision = $\max(U, \hat{U});$ end for

 $\hat{U} = \mathbb{E}\left(\hat{G}\right) \cdot e^{-\eta \frac{\hat{\beta}}{\zeta}}$

Receive

Get tasks, parameters and constraints

Update Update the demand and load

Estimate

Get the expected ranking and utilities for Decisions 1 & 2

Decide

University of Glasgow

Get the appropriate decision



he Scottish Informatics



Research Axis C TASKS





Challenges



Tasks Characteristics

- Load, constraints, processing
- \checkmark Define the contextual vector of tasks



Peers Characteristics

- ✓ Load, data, processing capabilities
- Define the Peer Contextual Vector (PCV): load, data relevance, speed of processing, communication cost



Efficiency of Allocation

✓ How can we match tasks contextual vector with PCVs



Uncertainty Management

- ✓ Select the appropriate technology
- ✓ Fuzzy Logic (FL) seems the solution

Contribution

- ✓ We define a function h(t_j, PCV_k) (k is the index of a peer node) that delivers a 'judgement' of the efficiency for the allocation
 - ✓ We realize h() with a Type-2 FL System (T2FLS)
- We handle the uncertainty in two axes: (i) in the definition of fuzzy sets; (ii) In the definition of membership functions
- We define the new concept of Type2D Sets, i.e., membership functions are automativcally defined by ML

Uncertainty Management



Type-2 Fuzzy Sets

Type-2 fuzzy sets and systems generalize standard Type-1 fuzzy sets and systems so that more uncertainty can be handled

University of Glasgow

he Scottish Informatics

Uncertainty Management





Reward and Decision Making





Algorithm

he Scottish Informatics

```
while true do
    if training interval expiration is true then
        trainT2DFLS();
    end
    if contextual vectors are reported by peers then
        updatePCVs();
    end
    if t_j is received then
        \mathbf{t}_i = \text{defineTaskVector()};
        for k \leftarrow 1 to N do
            PoA_k = \text{getT2DFLSResult}(l_k, s_k, \lambda_j);
            Z_k = \text{getReward}(PoA_k, r_k, \kappa_k);
            Z.add(Z_k);
        end
        sort(Z);
        select the best node and allocate t_j;
    end
end
```



Future Research Directions



he Scottish Informatics & Computer Science Alliance

THANK YOU

More Publications, Datasets, Presentations can be found at: http://kostasks.users.uth.gr http://www.iprism.eu Email: kostasks@uth.gr