



## **Résumé :**

L'Edge Computing (EC) fournit des ressources de stockage et de calcul distribuées à proximité des équipements IoT et mobiles afin d'offrir des services beaucoup plus interactifs en utilisant les capacités des réseaux de nouvelles générations. C'est une architecture prometteuse pour l'hébergement des services exigeants en termes de latence et de capacité qui sont difficiles à fournir dans les architectures cloud classiques. Par conséquent, les nouvelles architectures Edge-Cloud Mobiles (ECM) impliquant l'EC et le cloud computing émergent, vu leurs multiples avantages permettant d'obtenir un débit plus important, une latence faible, une richesse en termes de ressources, une plus grande réactivité et une évolutivité transparente, etc. Avec la miniaturisation des équipements, le développement rapide et l'utilisation largement répandue des réseaux sans fil, les appareils mobiles et IoT sont de plus en plus connectés et deviennent essentiels pour nos communications quotidiennes, nos activités personnelles et professionnelles. Toutefois, leurs ressources limitées exigent le recours au déchargement de calcul afin de supporter l'exécution multitâche d'applications utilisateur gourmandes en termes de ressources de calcul et de stockage. Pour ces raisons, l'edge-cloud mobile offre des services virtualisés de déchargement afin de libérer ces équipements des traitements lourds qui agissent au détriment de leurs ressources critiques. Cependant, ces services sont très sensibles à la latence et requièrent une gestion efficace des ressources, notamment de communication, de stockage et de calcul. Aussi, avec la mobilité des utilisateurs, le positionnement et le maintien de la qualité et de la continuité de ces services restent un vrai challenge.

Dans cette thèse, nous nous concentrons sur l'amélioration des services de déchargement de calcul dans des architectures EC et ECM composées de plusieurs serveurs. Notre objectif est de mettre en place des solutions algorithmiques efficaces permettant de prendre les meilleures décisions de déchargement des tâches, tout en gérant efficacement les problèmes d'allocation des ressources qui en résultent, en garantissant les contraintes de temps d'exécution des tâches mises en jeu et en minimisant la consommation énergétique. Toutes nos contributions considèrent des systèmes multiutilisateurs multiserveurs. Notre première contribution est une solution de déchargement distribuée qui considère des utilisateurs avec des tâches provenant d'une application déchargeable contrainte à un délai d'exécution strict. Cette exécution doit avoir lieu localement ou aux serveurs distants avec des décisions qui minimisent la consommation énergétique du système étudié, tout en considérant les contraintes des traitements et des ressources. Nous proposons dans notre deuxième contribution une extension de cette solution en introduisant l'indépendance des délais d'exécution des tâches de calcul qui, cette fois, proviennent de plusieurs applications locales. En plus, le nouveau modèle propose un déchargement adaptatif qui, en plus des contraintes des délais, prend en considération les besoins des dispositifs mobiles en termes de capacités de calcul locales. Après l'évaluation des solutions et l'analyse des résultats, notre troisième contribution étudie le problème de déploiement et de migration des services, car le positionnement de leurs conteneurs est très critique d'après les résultats de la deuxième contribution. Ainsi, nous proposons une première solution qui décide le déploiement et la migration des services virtualisés, tout en satisfaisant au maximum les exigences des utilisateurs liées à la bande passante de service. La deuxième solution de cette contribution propose une extension du modèle proposé qui introduit la pénalisation des conteneurs avec priorité afin de garantir l'obtention de solutions même lors des situations de ressources critiques. Cette solution détermine les meilleures décisions de placement de conteneurs qui minimisent une fonction de coût et de pénalisation, tout en favorisant les conteneurs prioritaires. Enfin, les problèmes de migration/déploiement obtenus nécessitent une optimisation combinatoire. À mesure que ces problèmes de décision sont d'une grande complexité, nous faisons appel à des méthodes heuristiques de complexités acceptables qui approchent efficacement leurs solutions.

**Mots clés :** Edge-Cloud Mobile, Cloud Computing, Edge Computing, Déchargement des Tâches, Services Virtualisés, Conteneurs, Migration, Déploiement, Recuit Simulé, Système de colonie de fourmis.

# OPTIMIZATION OF VIRTUALIZED SERVICES AND MULTIUSER MULTITASK OFFLOADING IN MULTISERVER MOBILE EDGE-CLOUD ENVIRONMENTS

## Abstract :

Edge computing provides distributed storage and computing resources near IoT and mobile devices to deliver much more interactive services using the capabilities of next-generation networks. It is a promising architecture for hosting services and which are demanding in terms of latency and capacity that are difficult to provide in traditional cloud architectures. As a result, new mobile edge-cloud architectures involving edge and cloud computing are emerging due to their multiple advantages. These last allow higher data-rates, low latency, resources richness, important responsiveness and seamless scalability, etc. With the devices miniaturization, rapid development and widespread use of wireless networks, mobile and IoT devices are increasingly connected and become essential for our daily communications, our personal and professional activities. However, their limited resources require the use of computation offloading to support the multitask execution of resource-hungry user applications regarding computations and storage. For these reasons, mobile edge-cloud offers virtualized offloading services to free these devices from heavy processing which negatively impacts their critical resources. However, these services are highly latency-sensitive and require efficient resource management including communication, storage, and processing. Also, with users mobility, positioning and maintaining the quality and the continuity of these services remains a real challenge.

In this thesis, we focus on improving computation offloading services in multi-servers edge and edge-cloud computing architectures. Our goal is to establish efficient algorithmic solutions allowing to make the best decisions to offload tasks. This is done while effectively allocating resources, ensuring tasks' execution time constraints, and minimizing energy consumption. All our contributions consider multiuser, multiserver systems. Our first contribution is a distributed offloading solution that considers users with many tasks generated by one offloadable application that is constrained to a strict execution deadline. Its execution must take place locally or at remote servers while the decisions minimize the energy consumption of the studied system and consider the processing and resources constraints. In our second contribution, we propose an extension of this first solution by introducing the tasks' execution times independence where these tasks can be generated from several local applications. In addition, the new model proposes an adaptive offloading approach which, besides time constraints, considers the needs of mobile devices in terms of local computing capacities. After evaluating the solutions and analyzing the results, the positioning of the containers' offloading services is very critical. Based on this finding, our third contribution addresses the services' deployment-migration problem. So we offer a first solution that decides on the deployment and the migration of virtualized services, while satisfying the maximum user requirements related to service bandwidth. The second solution of this contribution proposes an extension of the proposed model which introduces the penalization of containers with priority. This extension ensures solutions even during situations of critical resources. This new model determines the best placement decisions for containers that minimizes a cost and penalty function such that priority containers are favored. Finally, the obtained migration-deployment problems require combinatorial optimization. As these decision problems are too complex, we use heuristic methods of acceptable complexities that provide effective approximate solutions.

**Key Words:** Mobile Edge-Cloud, Cloud Computing, Edge Computing, Task Offloading, Virtualized Services, Containers, Migration, Deployment, Simulated Annealing, Ant Colony System.