



HAL
open science

Apprentissage en ligne par renforcement centré utilisateur pour la composition émergente d'applications ambiantes

Walid Younes, Françoise Adreit, Sylvie Trouilhet, Jean-Paul Arcangeli

► To cite this version:

Walid Younes, Françoise Adreit, Sylvie Trouilhet, Jean-Paul Arcangeli. Apprentissage en ligne par renforcement centré utilisateur pour la composition émergente d'applications ambiantes. Conférence Nationale d'Intelligence Artificielle en collaboration avec JSAI 2019 (CNIA 2019), Jul 2019, Toulouse, France. pp.179-186. hal-02879726

HAL Id: hal-02879726

<https://hal.science/hal-02879726v1>

Submitted on 24 Jun 2020

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Open Archive Toulouse Archive Ouverte

OATAO is an open access repository that collects the work of Toulouse researchers and makes it freely available over the web where possible

This is an author's version published in:

<http://oatao.univ-toulouse.fr/26345>

Official URL

https://afia.asso.fr/wp-content/uploads/2019/12/CNIA_2019.pdf

To cite this version: Younes, Walid and Adreit, Françoise and Trouilhet, Sylvie and Arcangeli, Jean-Paul *Apprentissage en ligne par renforcement centré utilisateur pour la composition émergente d'applications ambiantes*. (2019) In: Conférence Nationale d'Intelligence Artificielle en collaboration avec JSAI 2019 (CNIA 2019), 1 July 2019 - 5 July 2019 (Toulouse, France).

Any correspondence concerning this service should be sent to the repository administrator: tech-oatao@listes-diff.inp-toulouse.fr

Apprentissage en ligne par renforcement centré utilisateur pour la composition émergente d'applications ambiantes

Walid Younes¹

Françoise Adreit²

Sylvie Trouilhet¹

Jean-Paul Arcangeli¹

¹ Université de Toulouse, IRIT, Université Paul Sabatier

² Université de Toulouse, IRIT, Université Toulouse Jean Jaurès

Prénom.Nom@irit.fr

Résumé

Les systèmes cyber-physiques connectés et ambiants entourent l'utilisateur humain de services plus ou moins complexes, qui doivent être le plus possible adaptés à ses préférences et à la situation courante. Nous proposons de les construire automatiquement et à la volée par composition de services plus élémentaires présents dans l'environnement, sans expression préalable des besoins de l'utilisateur. Pour cela, avec la forte variabilité dynamique de l'environnement ambiant et des besoins, ce dernier doit être sollicité mais à minima. Afin de produire les connaissances nécessaires à la composition automatique et en l'absence de données initiales, nous proposons une solution d'apprentissage en ligne par renforcement à horizon infini qui apprend incrémentalement de et pour l'utilisateur. Elle est décentralisée au sein d'un système multi-agent chargé de l'administration et de la composition des services.

Mots-clés

Apprentissage en ligne, apprentissage par renforcement, feedback utilisateur, services logiciels, composition de services, système multi-agent, intelligence ambiante

Abstract

Cyber-physical and ambient systems surround the human user with services which are more or less complex, and must be as tailored as possible to her/his preferences and the current situation. We propose to build them automatically and on the fly by composition of more elementary services present in the environment, without prior expression of the user's needs. In a context of high dynamic variability of both the ambient environment and the needs, the user must be involved at the minimum. In order to produce the knowledge necessary for automatic composition in the absence of an initial guideline, we have designed a generic solution based on online reinforcement learning. It is decentralized within a multi-agent system in charge of the administration and the composition of the services, which learns incrementally from and for the user.

Keywords

Online learning, reinforcement learning, user feedback, service discovery, selection and composition, multi-agent system, ambient intelligence

1 Introduction

Les systèmes ambiants et mobiles sont composés d'appareils fixes ou mobiles reliés par un ou plusieurs réseaux de communication. Ces appareils hébergent des composants logiciels qui fournissent des services et qui peuvent, eux-mêmes, requérir d'autres services. Ces composants sont des briques logicielles qui peuvent être assemblées en connectant des services requis à des services fournis pour composer des applications plus complexes [20]. Par exemple, l'assemblage d'un composant d'interaction non dédié présent dans un smartphone (*e.g.*, un *curseur*, un *bouton* ou un composant de reconnaissance vocale), d'un *adaptateur* et d'une *lampe connectée* peut permettre de réaliser une application permettant à un utilisateur de contrôler l'éclairage ambiant.

Les composants (matériels et logiciels) sont en général multi-propriétaires et gérés de manière indépendante : ils sont développés, installés et activés indépendamment les uns des autres. En raison de la mobilité des appareils et des utilisateurs, ils peuvent apparaître ou disparaître selon une dynamique imprévisible, conférant aux systèmes ambiants et mobiles un caractère ouvert et instable. À cela s'ajoute le nombre souvent important de composants, source de difficultés lors du passage à l'échelle. Dans ce contexte, les assemblages de composants sont difficiles à concevoir, à entretenir et à adapter.

Plongé au sein de ces systèmes, l'utilisateur humain peut utiliser les services qui sont à sa disposition. L'intelligence ambiante [24, 7, 19] vise à lui offrir un environnement personnalisé, adapté à la situation, c'est-à-dire à lui fournir les bons services au bon moment, en anticipant ses besoins qui eux-mêmes peuvent évoluer. Pour cela, l'utilisateur peut être sollicité mais dans une limite raisonnable [3].

Notre projet a pour objectif de concevoir et de réaliser un "moteur de composition" qui détecte périodiquement les

composants présents dans l’environnement ambiant et qui les assemble dynamiquement et automatiquement afin de construire des applications “composites” adaptées à l’environnement et à l’utilisateur, c’est-à-dire opérationnelles, utiles et utilisables [9]. En l’absence de besoin explicité *a priori*, le moteur apprend les préférences de l’utilisateur en fonction de la situation. Notre approche rompt avec le traditionnel mode *top-down* pour le développement d’applications : la réalisation d’un assemblage n’est pas guidée par les besoins explicites de l’utilisateur ni par des plans d’assemblage prédéfinis ; au contraire, les applications composites sont construites à la volée en mode *bottom-up* à partir des services présents sur le moment dans l’environnement ambiant. Ainsi, les applications émergent de l’environnement, en tirant profit des opportunités. Dans ce cadre, l’utilisateur ne demande pas un service ou une application : les applications émergentes lui sont fournies en mode *push*¹. Pour cela, nous avons développé une solution basée sur l’apprentissage automatique. L’objectif de ce papier est de présenter les principes de cet apprentissage et de leur mise en œuvre. Il est organisé comme suit. L’architecture du système incluant le moteur de composition et l’utilisateur est présentée dans la section 2. Dans la section 3, notre problématique d’apprentissage est analysée en termes de motivations, d’objectifs et de données. Cette section se conclut par la définition du type d’apprentissage : en ligne et par renforcement. Les principes de notre solution d’apprentissage sont ensuite exposés dans la section 4. La section 5 résume l’état de l’art en matière d’apprentissage pour la composition logicielle automatique et positionne notre proposition. En conclusion, dans la section 6, le lecteur trouvera un bref résumé de la contribution, un point d’avancement sur le développement de la solution et une discussion sur les problèmes ouverts et la suite de ce travail.

2 Architecture du système de composition

Afin de répondre à l’exigence d’automatisation de la composition opportuniste de composants logiciels, nous avons défini une architecture logicielle du système de composition [25]. Cette section en présente les grandes lignes et la figure 1 en donne une vue simplifiée.

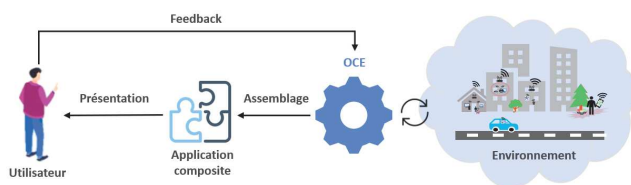


FIGURE 1 – Vue simplifiée de l’architecture

Au cœur du système, le moteur OCE (*Opportunistic Composition Engine*) a pour fonction de concevoir pour l’uti-

1. Un exemple de cas d’utilisation est développé dans [25].

lisateur des applications composites en assemblant des composants métiers et des composants d’interaction disponibles. Pour cela, il perçoit les composants et leurs services présents dans l’environnement ambiant, puis choisit d’établir des connexions entre des services requis et des services fournis faisant ainsi émerger les applications.

OCE propose à l’utilisateur les applications émergentes et c’est l’utilisateur qui décide *in fine* de leur déploiement : il peut accepter ou rejeter une application proposée ou encore la modifier [14]. Ce point est particulièrement important dans le domaine de l’interaction homme-machine pour lequel le contrôle par l’utilisateur de son environnement d’interaction est fondamental [3]. Après acceptation, l’application émergente est déployée automatiquement.

Afin de solliciter l’utilisateur *a minima*, le système de composition doit être le plus autonome possible. Notre solution se conforme aux principes de l’informatique autonome et au modèle MAPE-K (*Monitor, Analyse, Plan, Execute - Knowledge*) [12]. Cependant, alors que dans le modèle MAPE-K la phase d’exécution a pour fonction de réaliser ce qui a été planifié, l’application émergente est ici présentée à l’utilisateur sous la forme d’un assemblage de composants avant d’être éventuellement déployée. Les actions que l’utilisateur effectue en réponse sont sources de feedback et permettent à OCE d’apprendre et de faire évoluer ainsi ses connaissances (voir les sections suivantes).

Comme, fondamentalement, l’environnement ambiant est distribué et les composants matériels et logiciels sont administrés par différentes autorités, nous avons choisi de concevoir le moteur sous la forme d’un système multi-agent (SMA) dans lequel chaque service d’un composant est “administré” par un agent. De manière générale, les agents sont des entités autonomes qui coopèrent afin d’atteindre un but commun [11]. Ceux du moteur OCE ont pour rôle d’établir des connexions (ou des déconnexions) entre les services qu’ils administrent. Chaque agent a pour but de trouver un service partenaire avec qui se connecter. Pour cela, les agents échangent des messages dans le cadre du protocole d’interaction ARSA (*Advertise, Reply, Select, Agree*) [25] qui supporte la coopération entre agents en mode asynchrone dans un contexte dynamique et ouvert : ainsi, les agents continuent de coopérer même si un agent disparaît (en cas de désactivation du service administré ou de panne) ou si un message n’arrive pas à son destinataire.

3 Apprendre : pourquoi, quoi, à partir de quoi et comment ?

La fonction d’OCE est de concevoir des assemblages de composants logiciels qui réalisent des applications émergentes pertinentes, puis de les proposer à l’utilisateur et enfin de les déployer suivant les retours de l’utilisateur. Pour cela, le moteur doit prendre des décisions de connexion entre services dans l’objectif de satisfaire l’utilisateur.

Quatre éléments principaux contraignent ces prises de décision : la dynamique et l’imprévisibilité de l’environne-

ment ambiant, le nombre de composants et de services potentiellement important, la dynamique et la variabilité des besoins de l'utilisateur et enfin, la nécessité de solliciter *a minima* l'utilisateur.

Pour prendre des décisions pertinentes, le moteur OCE a besoin de connaissances qui, dans ce contexte, doivent être apprises automatiquement. Dans cette section, nous analysons pourquoi OCE doit apprendre, ce qu'il doit apprendre et quelles sont les données d'apprentissage.

3.1 Pourquoi apprendre ?

De par la dynamique et l'imprévisibilité de l'environnement ambiant, la combinatoire générée par le nombre important de composants ainsi que la dynamique et la variabilité de ses propres besoins, l'utilisateur n'est pas en mesure d'exprimer *a priori*, explicitement et d'une manière exhaustive ses besoins et ses préférences ni de les traduire en plans d'assemblages dans les différentes situations qu'il peut rencontrer. Par conséquent, OCE ne peut se baser ni sur des besoins de l'utilisateur explicités *a priori* ni sur des règles d'assemblage prédéfinies. OCE doit donc apprendre et ce, à partir de l'expérience.

3.2 Apprendre quoi ?

OCE doit construire les connaissances nécessaires pour proposer des applications pertinentes à l'utilisateur. Il doit apprendre ce que celui-ci préfère quand certains composants sont présents dans l'environnement ambiant. Basées sur ces préférences à un moment donné, ces connaissances seront exploitées pour prendre de futures décisions dans des situations identiques ou similaires. Par exemple, OCE peut apprendre que l'utilisateur préfère contrôler l'éclairage ambiant au moyen du curseur embarqué sur son smartphone plutôt qu'avec l'interrupteur mural connecté.

3.3 À partir de quoi apprendre ?

Si le moteur ne peut pas disposer initialement d'un ensemble de données d'apprentissage, en revanche, des données peuvent être observées au fil de l'utilisation du système pour servir de base aux décisions futures.

Plusieurs sources de données sont possibles. La première provient de l'observation et de l'exploitation des retours de l'utilisateur sur l'application émergente proposée (acceptation, modification ou rejet). Par exemple, en cas de modification par l'utilisateur du contrôle de l'éclairage ambiant, OCE peut modifier ses préférences de connexion entre services (le curseur du smartphone plutôt que l'interrupteur). On peut récupérer d'autres feedbacks à partir de l'éditeur utilisé par l'utilisateur pour manipuler les composants et leurs connexions [14], notamment en observant certaines actions de l'utilisateur (par exemple, pousser un composant en dehors de la fenêtre d'édition). En interne au SMA, il est également envisageable d'extraire des informations de feedback à partir des interactions des agents (par exemple, en cas d'échanges soutenus entre deux agents, on pourrait favoriser la connexion entre leurs services respectifs).

À ce stade de notre travail, nous avons choisi d'exploiter la première source de données qui traduit les préférences de l'utilisateur en fonction des composants présents dans l'environnement ambiant. En effet, nous faisons l'hypothèse que même si l'utilisateur ne peut pas expliciter *a priori* ses besoins, il est capable de réagir sur le moment à la proposition d'une application construite automatiquement. Il est alors possible de capturer cette réaction sous la forme de feedback et d'en extraire de la connaissance utile à des prises de décision futures. Ainsi, le moteur apprend à partir des interactions avec l'utilisateur sans le surcharger.

3.4 Comment apprendre ?

L'absence de données initiales et de solutions connues rend impossible un apprentissage supervisé ou non supervisé. De plus, la dynamique de l'environnement, avec les services qui apparaissent et disparaissent de manière imprévisible, rend très difficile voire impossible la construction d'un modèle statique de prédiction ou de classification. Pour ces raisons, nous proposons un apprentissage par adaptation progressive permettant au moteur OCE d'apprendre en continu en exploitant les informations de feedback que l'utilisateur fournit itérativement. C'est un apprentissage *en ligne, par renforcement et à horizon infini*. Ce type de solution permet à un système de s'adapter sur le long terme en interagissant avec son environnement [4].

L'objet de l'apprentissage est ici de contribuer à déterminer une action. Les agents d'OCE raisonnent en s'appuyant sur des connaissances construites et mises à jour en ligne, de manière incrémentale, au fur et à mesure de l'expérience et au gré des interactions avec l'utilisateur et des éventuelles évolutions de ses préférences. Selon le modèle de l'apprentissage en ligne [8], les agents font une "prédiction" (l'assemblage) et l'environnement (ici l'utilisateur) apporte une réponse. Cependant, le retour donné ici par l'utilisateur n'a pas le caractère d'exactitude de la réponse de l'environnement du modèle classique de l'apprentissage en ligne. Pour cette raison, nous hybridons les principes de l'apprentissage en ligne avec ceux de l'apprentissage par renforcement [8] : la réponse de l'utilisateur permet de renforcer les connaissances des agents, et les décisions à l'itération t s'appuient sur les connaissances cumulées lors des itérations précédentes. Enfin, considérant la dynamique, à la fois de l'environnement ambiant et de l'utilisateur, cet apprentissage doit être en plus à horizon infini [8], ce qui n'exclut pas des phases de stabilisation des connaissances.

Notons en complément que cette approche n'exclut pas non plus la possibilité d'exploiter des connaissances connues *a priori* (par exemple, des règles générales d'assemblage de composants métiers, des règles d'ergonomie d'assemblage de composants d'interaction) qui pourraient être fournies initialement et ainsi accélérer le processus d'acquisition des connaissances.

4 Principes de la solution

OCE est un système multi-agent dans lequel chaque agent administre un service d'un composant. Dans le cadre de l'architecture présentée en section 2, OCE opère dans un cycle, appelé *cycle moteur*. Les agents eux-mêmes fonctionnent de manière cyclique : dans un cycle classique, appelé *cycle agent*, un agent perçoit, puis décide et enfin agit. Plusieurs cycles agent s'effectuent ainsi dans un cycle moteur : les agents perçoivent les messages reçus, décident et effectuent les actions conformément au protocole ARSA. En fin de cycle moteur, après présentation d'un assemblage à l'utilisateur, OCE récupère le feedback de l'utilisateur pour apprendre. Il peut ensuite sonder une nouvelle fois l'environnement ambiant et démarrer un nouveau cycle avec les connaissances des agents mises à jour. Dans un premier temps, nous présentons ce que fait un agent dans un cycle agent, en particulier comment il exploite ses connaissances. Puis, nous expliquons comment chaque agent construit et fait évoluer ses connaissances par apprentissage, et nous discutons ce mode d'apprentissage.

4.1 Cycle de vie d'un agent : de la perception à l'action

Pour décider de l'action à effectuer, un agent construit une représentation de la situation courante dans la phase de perception. Il compare ensuite cette situation à des situations de référence qu'il a déjà rencontrées pour pouvoir l'évaluer et choisir l'agent à qui il répond (et donc l'action) dans la phase de décision. Enfin, il effectue l'action.

Construction de la situation courante. On appelle *situation courante* S_t^i pour un agent A^i la situation dans laquelle se trouve A^i dans le cycle moteur courant : elle est composée de l'ensemble des agents services perçus par A^i dans l'environnement et compatibles avec lui du point de vue de la composition². Cette situation est construite par A^i , incrémentalement dans un cycle moteur, à partir des messages qu'il reçoit. C'est à l'étape de perception qu'elle est créée puis actualisée conformément à l'algorithme 1.

La situation courante S_t^i est formée d'un ensemble des couples de la forme $(A^j, Type_Message)$ où :

- A^j est l'identifiant de l'agent émetteur du message,
- $Type_Message$ est le type de message envoyé dans le cadre du protocole ARSA c'est-à-dire *Advise*, *Reply*, *Select* ou *Agree*.

La situation courante est ensuite rapprochée des situations de référence qu'a déjà rencontrées l'agent.

Rapprochement avec les situations de référence. On appelle *situation de référence* pour un agent A^i une situation identifiée lors d'un précédent cycle moteur. Un agent dispose ainsi d'un ensemble de situations de référence qui constituent sa mémoire, noté Ref^i . Nous verrons plus loin comment cette connaissance est construite et maintenue par apprentissage. A l'instar d'une situation

2. Deux services sont dits compatibles si l'un est fourni S_F et l'autre est requis S_R , et si S_F rend le service S_R ($S_R \subset S_F$).

Algorithme 1 Perception d'un agent A^i

- 1: Récupérer le lot de messages reçus
 - 2: **pour** chaque message reçu **faire**
 - 3: Extraire le couple $(A^j, Type_Message)$
 - 4: **si** service(A^i) et service(A^j) sont compatibles **alors**
 - 5: Actualiser la situation courante : $S_t^i \leftarrow S_t^i \cup \{(A^j, Type_Message)\}$
 - 6: **fin si**
 - 7: **fin pour**
-

courante, une situation de référence est composée d'un ensemble d'agents services perçu par l'agent dans l'environnement à un moment donné, agents compatibles avec lui du point de vue de la composition

Une situation de référence de A^i est formée d'un ensemble de couples de la forme $(A^j, Score_j^i)$, où :

- A^j est l'identifiant de l'agent émetteur du message,
- $Score_j^i$ est une valeur numérique qui représente pour A^i l'intérêt d'une connexion avec le service administré par A^j .

Le rapprochement a pour objectif d'identifier la situation courante parmi les situations de référence ou, à défaut, de sélectionner les situations de référence "similaires" à la situation courante. Il s'effectue sur la base des identifiants des agents présents dans les situations, en faisant abstraction des types de messages et des scores. Il est réalisé par la fonction *Calculer_Similarité* qui construit un sous-ensemble avec les situations de référence dont le degré de similarité avec S_t^i est supérieur à un seuil ξ ³. Le sous-ensemble est réduit à S_t^i si cette situation est connue (elle fait déjà partie des situations de référence); il est vide si aucune situation similaire n'a été trouvée.

Étant donné l'ensemble Ref^i des situations de référence de l'agent A^i et un seuil ξ , la fonction *Calculer_Similarité* est définie de l'ensemble des situations courantes pour A^i , Sit^i , vers un ensemble de couples composés d'une situation de référence et d'un degré de similarité (1) : à la situation courante rencontrée à l'instant t, S_t^i , est associé l'ensemble des couples formés par une situation de référence SR_k^i et son degré de similarité d_k avec S_t^i , dans le cas où d_k est supérieur ou égal au seuil ξ .

$$Calculer_Similarité : Sit^i \rightarrow \mathcal{P}(Ref^i \times \mathbb{R}) \quad (1)$$

$$S_t^i \mapsto \{(SR_k^i, d_k)\}_{k \leq |Ref^i|} \text{ et } d_k \geq \xi$$

La fonction *Calculer_Similarité* peut implanter un algorithme de *clustering* pour regrouper les situations de référence en classes. Ceci permet d'exécuter l'opération de rapprochement plus rapidement⁴ : il suffira à l'agent de comparer la situation courante aux représentantes dans chacune des classes.

3. Pour donner une idée, le degré de similarité entre deux situations peut être calculé sur la base la proportion d'agents en commun.

4. En particulier dans le cas des environnements ambiants où le nombre de situations de référence peut être important.

Marquage de la situation courante. La phase de marquage permet d'enrichir la situation courante à l'aide des situations de référence qui ont été sélectionnées dans la phase précédente, en attribuant un score aux agents de la situation courante. A l'issue de ce marquage, l'agent A^i pourra établir un classement préférentiel des agents A^j de la situation courante et en sélectionner un auquel répondre. Le marquage est réalisé par la fonction *Marquer_Situation* qui attribue les scores à partir des situations retournées par la fonction *Calculer_Similarité*. Si la situation courante S_t^i est connue, la fonction *Calculer_Similarité* a renvoyé la situation de référence correspondante et les valeurs des scores sont reproduites à l'identique par la fonction *Marquer_Situation*. Sinon, la fonction *Marquer_Situation* calcule le score $Score_j^i$ de chaque agent A^j de S_t^i ⁵. Si l'un des agents de la situation courante n'apparaît pas dans les situations de référence (c'est le cas lors de l'apparition d'un service nouveau dans l'environnement ambiant), un score arbitraire lui est attribué. Le choix de la valeur de ce score permet de favoriser la prise en compte de la nouveauté⁶. Ce choix appartient à l'agent A^i . Il gère ce choix et le fait varier en fonction de ce qu'il a appris sur la sensibilité à la nouveauté de l'utilisateur. Dans le cas où aucune situation de référence similaire n'a été trouvée, les scores sont initialisés avec une valeur arbitraire. Comme précédemment, le choix de cette valeur initiale appartient à l'agent A^i .

Choix de l'agent. Dans cette phase, l'agent A^i sélectionne un agent A^j de la situation courante (A^i répondra alors au message de A^j). Cette opération est effectuée par la fonction *Choisir_Agent* qui prend en paramètre la situation courante marquée et renvoie l'agent qui maximise un critère d'optimisation. Pour cela, plusieurs stratégies sont possibles basées sur les scores ou le type de message (dans l'ordre *Agree*, *Select*, *Reply*, *Advertise*) ou une combinaison de ces deux critères.

L'algorithme 2 synthétise le comportement de l'agent lors de l'étape de décision.

Algorithme 2 Décision d'un agent A^i

- 1: $Sit_Similaires_t^i \leftarrow$
 Calculer_Similarité(S_t^i)
 - 2: $Situation_Marquée_t^i \leftarrow$
 Marquer_Situation($S_t^i, Sit_Similaires_t^i$)
 - 3: $A^j \leftarrow$ *Choisir_Agent*($Situation_Marquée_t^i$)
-

Réalisation de l'action. Dans cette étape, l'agent A^i donne suite au message de l'agent A^j choisi à l'étape précédente. Le type de message à envoyer suit, dans le cadre du protocole ARSA, celui de A^j . Ainsi, un agent continue à coopérer avec les autres agents jusqu'au traitement d'un message de type *Agree*. Dans ce cas, il accepte la

5. Ce calcul peut être la moyenne des scores de A^j dans les situations de référence sélectionnées, pondérée par les degrés de similarité.

6. Par exemple, en choisissant une valeur supérieure aux scores des autres agents de la situation.

connexion avec le service géré par A^j et se met en attente d'un feedback sur sa décision. L'algorithme 3 décrit le comportement de A^i dans la phase d'action.

Algorithme 3 Comportement d'un agent service

- 1: **selon** (*Type_Message*)
 - 2: **cas** *Advertise* :
 - 3: envoyer un message *Reply* à A^j ;
 - 4: **cas** *Reply* :
 - 5: envoyer un message *Select* à A^j ;
 - 6: **cas** *Select* :
 - 7: envoyer un message *Agree* à A^j ;
 - 8: **cas** *Agree* :
 - 9: réaliser la connexion avec A^j ;
 - 10: se mettre en attente d'un feedback ;
 - 11: **fin selon**
-

4.2 Apprentissage par feedback utilisateur

Au cours des différents cycles agent (perception-décision-action), l'agent exploite ses connaissances. L'apprentissage de ces dernières s'effectue hors cycle agent, lorsque le cycle moteur arrive à son terme et que le feedback de l'utilisateur est renvoyé à OCE. Ce feedback porte sur l'ensemble de l'assemblage proposé. Il est répercuté sur les agents qui composent l'assemblage. Chacun de ces agents A^i construit alors une nouvelle situation de référence à partir de sa situation courante marquée. Pour cela, A^i calcule une valeur de renforcement du score de chaque agent A^k de la situation courante, notée $r_{A^k}^i$. Le calcul des $r_{A^k}^i$ utilise une variable $\beta > 0$, dont le choix appartient à A^i .

Trois cas sont possibles :

1. *L'utilisateur a accepté l'assemblage dans son intégralité.* La décision prise par tout agent A^i de l'assemblage de choisir son partenaire A^j est bonifiée :

$$r_{A^j}^i = \beta$$

$$r_{A^k}^i = 0, \forall A^k \in S_t^i \text{ t.q. } k \neq j.$$
2. *L'utilisateur a rejeté l'assemblage dans son intégralité.* La décision prise par tout agent A^i de l'assemblage de choisir A^j est pénalisée :

$$r_{A^j}^i = -\beta$$

$$r_{A^k}^i = 0, \forall A^k \in S_t^i \text{ t.q. } k \neq j.$$
3. *L'utilisateur a modifié l'assemblage.* Si, au moyen de l'éditeur à sa disposition (cf. section 3.3), il a déconnecté le service géré par un agent A^i du service géré par un agent A^j et reconnecté le premier au service géré par un agent A^h , on bonifie cette connexion et on pénalise celle proposée par OCE :

$$r_{A^h}^i = (Score_j^i - Score_h^i) + \beta$$

$$r_{A^j}^i = -\beta$$

$$r_{A^k}^i = 0, \forall A^k \in S_t^i \text{ t.q. } k \neq h \text{ et } k \neq j.$$

Les autres connexions de l'assemblage proposé sont traitées comme une acceptation (cas 1) si elles ne sont pas modifiées ou comme un rejet (cas 2) si elle sont supprimées.

Pour calculer le score $Score_k^i$ des agents A^k figurant dans la situation courante S_t^i , l'agent A^i utilise la formule (2), inspirée de celle des algorithmes de bandits, dans laquelle $Score_k^i$ est la valeur du score de l'agent A^k et $\alpha \in [0, 1]$ est le facteur d'apprentissage :

$$Score_k^i = Score_k^i + \alpha(r_{A^k}^i - Score_k^i) \quad (2)$$

Dans le contexte dynamique et d'imprévisibilité de notre problème, il est difficile pour un agent de se projeter sur les prochaines situations de connexion et donc de déterminer l'importance des récompenses futures. Aussi, à ce stade de notre travail, nous faisons l'hypothèse que l'action choisie n'a pas d'effet sur les récompenses futures. Dans le cas contraire, il faudrait revenir à la formule plus complète du *Q-Learning* avec un facteur d'actualisation non nul.

Dans la formule, $(1 - \alpha)Score_k^i$ représente la part d'information que l'agent A^i garde de son expérience passée et $\alpha r_{A^k}^i$ celle qu'il apprend dans le cycle courant. On peut noter que le score des agents non retenus dans l'assemblage diminue systématiquement. Pour l'agent sélectionné A^j , en fonction de la valeur de β , le score peut aussi diminuer mais dans une moindre proportion ; la position de A^j est donc renforcée par rapport aux agents non retenus.

Une fois la situation de référence construite, l'agent A^i la stocke dans sa base de connaissances, Ref^i . Dans le cas où cette situation est déjà dans Ref^i (cas où la situation courante correspondait à une situation déjà connue), A^i se contente de mettre à jour les scores.

4.3 Discussion

Le but de l'apprentissage d'OCE est de maximiser la satisfaction de l'utilisateur. Notre solution ne demande pas la définition explicite des critères de satisfaction : OCE se base sur le feedback de l'utilisateur (qui accepte, refuse ou modifie un assemblage) et non sur l'évaluation de critères de qualité prédéfinis. Ceci confère à OCE un caractère générique (quel que soit l'utilisateur) et évolutif (l'utilisateur peut évoluer et ses critères de satisfaction aussi).

L'apprentissage ne porte pas sur l'algorithme de décision d'un agent (hors la part d'exploration inhérente à l'apprentissage par renforcement, un agent cherche toujours à se connecter avec le "meilleur" agent). Il porte sur la construction et l'adaptation des connaissances qui amènent l'agent à améliorer ses propositions. Ainsi, l'assemblage proposé par OCE pour une situation déjà rencontrée pourra différer de celui qu'il avait proposé auparavant parce qu'entre temps il a appris certaines préférences de l'utilisateur.

Dans le processus de composition d'un assemblage, les données remontent des agents et l'assemblage émerge des propositions locales de connexions. Ce comportement est caractéristique des systèmes multi-agents, dans lesquels la fonction du système n'est explicitement définie dans aucun des agents mais émerge de leurs interactions. De la même façon, l'apprentissage est distribué entre les agents. C'est l'agent qui apprend localement en révisant ses connaissances sur ses connexions locales, en ajoutant ou en modi-

fiant des situations de référence voire en oubliant certaines situations (par exemple celles qui sont trop anciennes). Ces situations représentent la vision locale de l'agent sur les organisations possibles au sein du système multi-agent.

L'apprentissage est donc ici un apprentissage concurrent dans lequel chaque agent est un apprenant autonome influencé par l'environnement [4]. On peut toutefois se demander si la vision purement locale est suffisante. Les agents ne pourraient-ils pas coopérer davantage et échanger, par exemple, des situations de référence pour une cohérence plus forte de leurs décisions ? Dans le même ordre d'idée, les agents qui gèrent les services d'un même composant "hôte" gagneraient probablement à se coordonner (par exemple, il peut être inutile d'annoncer un service fourni par un composant tant que les services requis par le même composant hôte ne sont pas satisfaits). Ainsi, avec des agents qui apprendraient sur d'autres agents, on s'orienterait vers un apprentissage multi-agent [2]. Il faut cependant noter la contribution de l'utilisateur à la cohérence de la décision globale : il évalue et contrôle les décisions d'OCE et son feedback est distribué aux agents. Transformé en connaissance, ce feedback global cadre les décisions des agents et donne une cohérence globale à l'agrégation des décisions individuelles.

5 Travaux connexes

Fondamentalement, les composants logiciels et les services sont des unités logicielles développées et déployées pour être réutilisées et composées. L'automatisation de la composition est une problématique largement traitée dans la littérature, en particulier en ce qui concerne les services Web. Pour F. Morh [16], le problème de la composition automatique des services se divise en deux grandes classes selon que la "structure" de la composition est connue au préalable ou non. Dans le premier cas, il s'agit de trouver à l'exécution les différents services qui vont permettre de réaliser un modèle donné (plan de composition, *workflow* de services...) en l'adaptant au mieux à la situation. Par exemple, MUSIC [18] permet une adaptation contextuelle de la composition en se basant sur un modèle : les plans sont sélectionnés au moment de l'exécution afin de maximiser un critère de qualité. Dans le deuxième cas, on crée de nouveaux services qui satisfont des pré-conditions et des post-conditions, ou des besoins explicités en amont. Dans tous les cas, la composition s'effectue en mode *top-down* (à l'inverse de notre approche *bottom-up*), à partir de modèles prédéfinis ou de buts explicites.

Dans [21], les auteurs présentent un état de l'art sur la composition de services en intelligence ambiante. Les solutions présentées reposent sur la formulation (sous différentes formes) d'un but à atteindre. Dans certains cas, l'utilisateur peut être impliqué et choisir une solution de composition parmi différentes compositions possibles qui lui sont présentées, comme c'est le cas dans notre architecture. Il n'est cependant pas présenté de système de composition qui repose sur l'apprentissage et, par conséquent, de contri-

bution de l'utilisateur à un processus d'apprentissage.

Dans [17], les auteurs proposent une solution à base d'apprentissage pour l'adaptation en ligne et en continu de systèmes logiciels à composants. À partir d'un but explicite et d'un ensemble de configurations connues satisfaisant ce but, il s'agit de trouver la meilleure selon des critères extrafonctionnels et non de faire émerger une fonctionnalité comme nous le proposons. L'adaptation n'est pas programmée mais apprise par renforcement à partir d'expérimentations sur les différentes configurations possibles. L'utilisateur est sollicité pour expérimenter mais pas pour donner un feedback explicite. Pour cela, les applications sont instrumentées et c'est l'environnement d'exécution qui génère les données de feedback.

Les travaux basés sur l'apprentissage pour automatiser la composition prennent le plus souvent la qualité de service (*QoS*) comme critère (par ex., [13]). Dans ce cadre, Wang *et al.* proposent une composition auto-adaptative de services Web en environnement dynamique afin de maximiser la *QoS* globale de la composition offerte à l'utilisateur [23]. Pour cela, en s'inspirant de [6, 10], ils modélisent la structure d'une composition de service sous la forme d'un processus de décision markovien contenant plusieurs *workflows*, la politique optimale pour le choix du meilleur *workflow* étant basée sur un algorithme de *Q-Learning*. La *QoS* n'est actuellement pas prise en compte dans notre solution ; ce pourrait être un élément de nature à enrichir la décision, obtenu par observation ou *via* le feedback de l'utilisateur.

Pour traiter l'évolutivité, le passage à l'échelle et l'optimisation dynamique de la composition, Wang *et al.* [22] étendent le travail exposé dans [23] et proposent un *framework* multi-agent collaboratif où les agents apprennent par renforcement en utilisant l'algorithme *Q-Learning*. Ici, le partage d'expérience entre les agents améliore l'efficacité et la vitesse d'apprentissage. Dans [15], les auteurs proposent une approche collaborative basée sur un algorithme d'apprentissage par renforcement appelé "Learning automata", pour adapter la composition de services Web et maintenir une *QoS* satisfaisante de la composition. Y. Charif et N. Sabouret utilisent eux aussi un protocole de coordination entre agents pour la chorégraphie dynamique de services : un agent dialogue à l'aide de requêtes et utilise un historique des conversations en guise de mémoire [5]. Ces approches coopératives sont intéressantes et rejoignent nos perspectives discutées en section 4.3.

6 Conclusion

Pour construire des applications ambiantes par assemblage automatique et dynamique de services logiciels, notre approche se démarque des solutions existantes en faisant émerger les applications de l'environnement ambiant en mode *bottom-up*, sans expression préalable des besoins ni modèle d'assemblage prédéfini. Afin de satisfaire au mieux l'utilisateur tout en limitant sa contribution dans un contexte de forte dynamique et d'imprévisibilité, cette émergence doit être contrôlée. Pour cela, notre "moteur de

composition" apprend en ligne, par renforcement et à horizon infini. C'est un système multi-agent dans lequel chaque agent gère un service. Les agents interagissent et coopèrent dans le cadre du protocole ARSA. Ils apprennent individuellement et en concurrence : chacun construit, à partir des messages qu'il reçoit, sa représentation locale de l'état du monde qui l'entoure. Il identifie ainsi la situation courante et la rapproche de situations rencontrées par le passé, ces dernières ayant été évaluées à partir de données de feedback de l'utilisateur. Le rapprochement entre situations permet à l'agent de prendre des décisions de connexion pertinentes pour le service qu'il administre. Cette méthode se rapproche du raisonnement par cas [1] basé sur la réutilisation de cas antérieurs et de solutions à un problème semblable au problème courant, apprises précédemment. En outre, le protocole ARSA et le rapprochement de situations permettent d'intégrer des services nouveaux et inconnus qui apparaissent soudainement dans l'environnement.

L'architecture fonctionnelle de notre système et le protocole ARSA ont été implémentés. Un prototype d'éditeur pour l'interaction avec l'utilisateur a été développé par ailleurs [14]. Avant de mettre en œuvre la solution d'apprentissage, il convient de finaliser la définition des différentes fonctions (calcul de similarité, marquage de situation, choix de l'agent) et les paramètres d'apprentissage. Nous pourrions alors conduire différentes expérimentations, portant sur différents cas d'utilisation identifiés, afin d'évaluer et de calibrer la solution.

Quelques questions restent ouvertes. D'une part, le temps d'apprentissage nécessaire pour que le moteur propose à l'utilisateur des applications utiles et utilisables pourrait s'avérer important. Une solution brièvement exposée en fin de section 3.4 pourrait consister à incorporer des règles métier ou des règles d'ergonomie au niveau des agents, sous la forme de situations de référence prédéfinies, afin d'accélérer l'acquisition des connaissances. Par ailleurs, notre solution d'apprentissage se base sur un feedback utilisateur capturé en phase de présentation de l'application. D'autres sources de données d'apprentissage sont possibles (cf. section 3.3). Dans tous les cas, un équilibre doit être respecté entre, d'un côté, la qualité et la quantité du feedback et de l'autre, la nature et la fréquence de la sollicitation de l'utilisateur. La précision de la valuation des situations à partir du feedback et le rapprochement de situations sont d'autres points critiques pour la qualité de la décision.

D'autre part, la formule (2), inspirée de celle du *Q-Learning*, ne prend pas en compte le facteur d'actualisation γ (cf. section 4.2). Cette question doit être étudiée plus précisément sur la base de cas concrets afin de déterminer en quoi le choix d'une "action" pourrait influencer sur les futures récompenses. Si tel était le cas, nous devrions réintroduire le facteur d'actualisation.

Enfin, comme nous l'avons indiqué en section 4.3, notre apprentissage est essentiellement individuel, sans échange d'information ni coordination entre les agents apprenants. L'introduction dans le système de mécanismes adéquats

devrait aussi améliorer la qualité de l'apprentissage, donc de la décision au niveau de l'ensemble des agents, et par conséquent la qualité de l'assemblage global.

7 Remerciements

Ce travail est en partie financé par la région Occitanie et le programme opérationnel FEDER-FSE Midi-Pyrénées et Garonne ainsi que par l'Université Paul Sabatier dans le cadre de l'opération neOCampus.

Références

- [1] A. Aamodt and E. Plaza. Case-based reasoning : Foundational issues, methodological variations, and system approaches. *AI Communications*, 7(1) :39–59, March 1994.
- [2] S. Albrecht and P. Stone. Autonomous Agents Modelling Other Agents : A Comprehensive Survey and Open Problems. *Artificial Intelligence*, 258 :66–95, 2018.
- [3] C. Bach and D. Scapin. Adaptation of ergonomic criteria to human-virtual environments interactions. In *Proc. of Interact'03*, pages 880–883. IOS Press, 2003.
- [4] J. Boes. *Apprentissage du contrôle de systèmes complexes par l'auto-organisation coopérative d'un système multi-agent : application à la calibration de moteurs à combustion*. PhD thesis, Université de Toulouse, UPS, 2014.
- [5] Y. Charif and N. Sabouret. Un protocole de coordination d'agents introspectifs pour la chorégraphie dynamique de services. *Revue d'Intelligence Artificielle (RSTI-RIA)*, 23(1) :47–79, 2009.
- [6] K. Chen, J. Xu, and S. Reiff-Marganiec. Markov-HTN Planning Approach to Enhance Flexibility of Automatic Web Service Composition. In *IEEE Int. Conf. on Web Services*, pages 9–16. IEEE, 2009.
- [7] D. J. Cook, J. C. Augusto, and V. R. Jakkula. Ambient intelligence : Technologies, applications, and opportunities. *Pervasive and Mobile Computing*, 5(4) :277–298, 2009.
- [8] A. Cornuéjols, L. Miclet, and V. Barra. *Apprentissage artificiel : Deep learning, concepts et algorithmes*. Eyrolles, 3ème édition, 2018.
- [9] J. Coutaz. Essai sans prétention sur l'Interaction Homme-Machine et son évolution. *1024 : Bulletin de la Société Informatique de France*, (1) :15–33, 2013.
- [10] P. Doshi, R. Goodwin, R. Akkiraju, and K. Verma. Dynamic workflow composition : Using markov decision processes. *Int. Journal of Web Services Research (IJWSR)*, 2(1) :1–17, 2005.
- [11] J. Ferber. *Multi-agent systems : An introduction to distributed artificial intelligence*. Addison Wesley, 1999.
- [12] J. O. Kephart and D. M. Chess. The vision of autonomous computing. *Computer*, 36(1) :41–50, Jan. 2003.
- [13] M. E. Khanouche, F. Attal, Y. Amirat, A. Chibani, and M. Kerkar. Clustering-based and QoS-aware services composition algorithm for ambient intelligence. *Information Sciences*, 482 :419–439, 2019.
- [14] M. Koussaifi, S. Trouilhet, J.-P. Arcangeli, and J.-M. Bruel. Ambient intelligence users in the loop : Towards a model-driven approach. In *Software Technologies : Applications and Foundations*, pages 558–572. Springer, 2018.
- [15] G. Li, D. Song, L. Liao, F. Sun, and J. Du. Learning automata-based adaptive web services composition. In *5th IEEE Int. Conf. on Software Engineering and Service Science (ICSESS)*, pages 792–795. IEEE, 2014.
- [16] F. Morh. *Automated Software and Service Composition*. SpringerBriefs in Computer Science. Springer, 2016.
- [17] R. Rodrigues Filho and B. Porter. Defining emergent software using continuous self-assembly, perception, and learning. *ACM Trans. on Autonomous and Adaptive Systems*, 12(3) :16 :1–16 :25, October 2017.
- [18] R. Rouvoy, P. Barone, Y. Ding, F. Eliassen, S. O. Hallsteinsen, J. Lorenzo, A. Mamelli, and U. Scholz. MUSIC : middleware support for self-adaptation in ubiquitous and service-oriented environments. In *Software Engineering for Self-Adaptive Systems*, volume 5525 of *LNCS*, pages 164–182. Springer, 2009.
- [19] F. Sadri. Ambient intelligence : A survey. *ACM Computing Surveys*, 43(4) :1–66, October 2011.
- [20] I. Sommerville. Component-based software engineering. In *Software Engineering*, chapter 16, pages 464–489. Pearson Education, 10th edition, 2016.
- [21] T. G. Stavropoulos, D. Vrakas, and I. Vlahavas. A survey of service composition in ambient intelligence environments. *Artificial Intelligence Review*, 40(3) :247–270, September 2011.
- [22] H. Wang, X. Wang, X. Hu, X. Zhang, and M. Gu. A multi-agent reinforcement learning approach to dynamic service composition. *Information Sciences*, 363 :96–119, 2016.
- [23] H. Wang, X. Zhou, X. Zhou, W. Liu, W. Li, and A. Bouguettaya. Adaptive service composition based on reinforcement learning. In *Proc. of the Int. Conf. on Service-Oriented Computing (ICSOC)*, pages 92–107. Springer, 2010.
- [24] M. Weiser. The computer for the 21st century. *Scientific American*, 265 :94–104, 1991.
- [25] W. Younes, S. Trouilhet, F. Adreit, and J.-P. Arcangeli. Towards an intelligent user-oriented middleware for opportunistic composition of services in ambient spaces. In *Proc. of the 5th Workshop on Middleware and Applications for the Internet of Things (M4IoT)*, pages 25–30, New York, NY, USA, 2018. ACM.