

Table des matières

1	Introduction générale	3
2	Assistants virtuels intelligents	4
2.1	Introduction	4
2.2	L'importance du contexte pour un SPA	6
2.3	Caractéristiques principales d'un SPA	6
2.3.1	Sensible au contexte	6
2.3.2	Évolutif	7
2.3.3	Multimodal	7
2.3.4	Anthropomorphe	7
2.3.5	Multi-plateforme et Flexible	8
2.4	Domaines d'applications des SPAs	8
2.4.1	Vie quotidienne	8
2.4.2	Assistance professionnelle	8
2.4.3	E-Apprentissage	9
2.5	Exemples de SPAs	9
2.5.1	Google assistant	9
2.5.2	Apple Siri	11
2.5.3	Amazon Alexa	12
2.5.4	Microsoft Cortana	12
2.6	Conclusion	12
3	Composants de base d'un assistant personnel	14
3.1	Introduction	14
3.2	Architecture d'un SPA	14
3.3	Reconnaissance automatique de la parole (ASR)	14
3.3.1	Définition	14
3.3.2	Modèle acoustique	14
3.3.3	Modèle de langage	14
3.3.4	Méthodes utilisées	14
3.4	Compréhension du langage naturel (NLU)	15
3.4.1	Définition	15
3.4.2	Classification d'Intent	15
3.4.3	Extraction d'entités	15
3.4.4	Analyse sémantique	15
3.5	Gestion du dialogue	15
3.5.1	Processus de décision Markovien (MDP)	16
3.5.2	État du gestionnaire de dialogue	17

3.5.3	Politique de gestion de dialogue	19
3.5.4	Gestion de dialogue par apprentissage	19
3.5.5	Apprentissage par renforcement	20
3.5.6	Simulateur d'utilisateur	20
3.6	Génération du langage naturel (NLG)	20
3.6.1	Détermination du contenu	21
3.6.2	Structuration de texte	21
3.6.3	Agrégation de phrases	22
3.6.4	Lexicalisation	22
3.6.5	Génération d'expressions référentielles (REG)	22
3.6.6	Réalisation linguistique	22
3.6.7	Systèmes basés encodeur-décodeur	24

To correct

To explain more deeply

Chapitre 1

Introduction générale

- Ici on parlera des motivations qui ont aboutis à ce projet, des objectifs de ce dernier ainsi que ses perspectives

Chapitre 2

Assistants virtuels intelligents

Introduction

Depuis la commercialisation du premier ordinateur grand public (Xerox PARC Alto) en 1973, le monde découvrit pour la première fois ce qui allait devenir l'apparence basique de chaque ordinateur moderne. En effet, la compagnie Xerox fut la première à proposer une interface graphique dotée de fenêtres, d'icônes et d'une souris pour se déplacer et d'un clavier pour écrire du texte. Bien que basique, cette idée lança alors plusieurs autres grandes marques sur le même chemin (IBM, Apple, Compaq ...). Par la suite, beaucoup ont essayé d'améliorer la façon dont l'homme utilisait sa machine : souris plus précise, écran doté d'une plus grande résolution, clavier plus enrichi, voire même l'introduction des écrans tactiles dans certains systèmes embarqués.

Cependant, certains voyaient encore cette façon d'utiliser la machine comme trop primitive, et peu intuitive. En effet laissez un enfant devant un ordinateur et il prendrait un bon moment pour apprendre à éditer ne serait ce qu'un simple fichier. Pour citer Donald A. Norman :

“We must design for the way people behave, not for how we would wish them to behave.”[1]

que nous pouvons traduire par :

“Nous devons concevoir selon le comportement des utilisateurs, et non pas selon la façon dont nous voudrions qu'ils se comportent.”

L'humanité a fait beaucoup de chemin depuis les années 70, l'utilisation d'un ordinateur de nos jours avec les moyens classiques (souris, clavier, écran ...) est devenue une tâche triviale, voire même une **seconde nature, cela reste cependant dû au fait que de plus en plus de jeunes enfants sont exposés depuis leur plus jeune âge au monde technologique qui les entoure, le processus d'apprentissage reste cependant présent, l'effort d'utiliser les outils communs reste lui aussi présent.**

La plus naturelle et plus ancienne façon de communiquer pour l'homme a toujours été la parole. Le développement de langues toutes aussi riches et complexes les unes que les autres a permis à l'humanité de briser plusieurs **barrières sociales**. L'avancement le plus naturel pour cette façon de communiquer serait donc de l'étendre aux machines que l'homme a su construire et améliorer au fil des années.

Motivé par cette manière que l'on a de communiquer entre nous, et épaulé par les récentes technologies telles que l'apprentissage automatique, le traitement automatique du langage naturel et

l'intelligence artificielle, les plus brillants des chercheurs ont entamé leurs travaux dans cette toute nouvelle direction.

Les Assistants Virtuels Intelligents (Smart Personal Assistant, SPA [2]) sont donc le produit de plusieurs années de recherche, visant tout d'abord à faciliter certaines tâches pour l'utilisateur. Les premiers SPAs étaient conçus comme des agents de conversation ou Chatbots, limités dans leurs actions et dépendant toujours d'un moyen de communication textuel, ce n'était pas la forme désirée du SPA. Avec l'avancement des recherches sur la reconnaissance automatique de la parole (Automatic Speech Recognition, ASR) et l'émergence de l'apprentissage automatique, les tout premiers assistants virtuels utilisant l'ASR étaient spécialisés dans certains domaines comme des systèmes médicaux d'aide à la décision. Il a ensuite été plus aisé de briser la barrière et de réaliser ce qui était encore une esquisse d'un SPA personnalisé. Aujourd'hui, et ce depuis l'avènement de l'apprentissage profond et la popularisation des Smartphones, de nouveaux SPAs comme Apple Siri (voir 2.5.2) et Google assistant (voir 2.5.1) et Amazon Alexa (voir 2.5.3) ont fait leurs apparitions, offrant de plus en plus de services personnalisés et spécifiques à chaque utilisateurs.

Dans la suite de ce chapitre nous essayerons de mieux détailler ce qu'est un SPA, ce qui est demandé d'un tel système, ses domaines d'application, en enchaînant par une description d'une pseudo-architecture potentielle de ce système, pour enfin conclure sur les limitations actuelles et les motivations de ce projet.

L'importance du contexte pour un SPA

Informellement, un SPA est un type d'agent (voir 2.3.2) logiciel qui peut effectuer certaines tâches et proposer des services dédiés aux utilisateurs qui vont d'une simple tâche (Ouvrir une fenêtre, lancer une application ...) à la réalisation de requêtes un peu plus complexes comme réserver une table dans un restaurant en passant un appel vocal (voir 2.5.1.1). Pour répondre efficacement à toutes sortes de requêtes, un SPA se doit donc de garder trace du contexte courant de sa conversation avec l'utilisateur. Il doit disposer d'un système capable d'enregistrer les informations pertinentes et de savoir les réutiliser, mais aussi de pouvoir déduire lesquelles de ces informations sont manquantes. On parle ici de Context-Awareness ou Sensibilité au contexte, comme vu dans [2].

D'après [2] et [3], *Day* et *Abwod* définissent un contexte comme suit :

“A context is any information that can be used to characterize the situation of an entity. An entity is a person, place, or object that is considered relevant to the interaction between a user and an application, including the user and applications themselves”

qui peut être traduit par :

“Un contexte est une information qui peut être utilisé pour caractériser l'état d'une entité. Une entité peut être une personne une place ou un objet, considérée comme pertinente à l'interaction entre l'utilisateur et l'application, ainsi qu'à ces deux derniers eux mêmes”

Il en découle que pour parvenir à développer un système qui puisse répondre aux besoins individuels et spécifiques de chaque personne, modéliser et prendre en compte le contexte semble être une solution prometteuse.

Caractéristiques principales d'un SPA

À partir de [2], nous pouvons dégager certaines caractéristiques principales qui peuvent être vues comme primordiales pour qualifier un assistant virtuel comme étant intelligent.

Sensible au contexte

Comme précédemment vu dans la définition du contexte (section 2.2), ce dernier peut être interprété comme tout aspect d'une entité (position d'un objet, couleur d'un objet, température d'une chambre, etc.). Un assistant dit intelligent doit donc être capable de capturer le concept du contexte, d'utiliser et de traiter toute information catégorisée comme contextuelle. Pour être plus précis, un SPA doit être sensible à l'évolution du contexte courant, par le biais de capteurs optiques, de microphones, ou tout ce qui pourrait amener l'utilisateur à faire évoluer la requête qu'il a émise. L'assistant devra donc proposer un système de mise à jour du contexte pour éliminer les informations inutiles et garder celles qui pourraient aider à répondre à la requête de l'utilisateur.

Évolutif

Comme vu dans la section 2.2, un SPA peut être vu comme un type d'agent. Pour rappel, d'après *Russel* et *Norvig* dans [4], un agent est une entité autonome pouvant interagir avec son environnement afin d'accomplir certaines tâches et peut être de plusieurs types :

- Agent à réflexes simples : agent exécutant ses actions à base de règles conditionnelles simples (c.à.d Si *Condition* alors *exécuter actions*), ils sont ainsi très simplistes et limités dans la portée de leurs actions.
- Agent basé modèle : semblable aux agents à réflexes simples, il est doté d'un modèle interne complexe censé représenter le monde extérieur auquel l'agent a accès. Cependant, il applique les actions de la même manière que le précédent type d'agents.
- Agent à but : ce type représente une amélioration des agents simples puisqu'il est doté d'un ensembles d'états buts à atteindre d'une façon ou d'une autre.
- Agent à utilité : il s'agit ici agents à buts qui tentent d'aboutir à leurs buts d'une manière optimisée (intelligente) utilisant une fonction de mesure adéquate pour le choix des différents états à atteindre.
- Agent apprenant : agent à utilité enrichi par un module d'apprentissage qui sert de juge pour répondre aux "critiques" des actions qu'il entreprend. Le terme agent évolutif est aussi employé.

Pour ce qui est des SPAs, les plus récents systèmes (ex : Amazon Alexa qui améliore son module de reconnaissance de la parole après chaque réponse non *réfutée* par l'utilisateur) peuvent être considérés comme des agents apprenants, répondant de ce fait à la contrainte évolutive imposée. Cependant, le domaine de l'auto-évolution des systèmes intelligents est encore un domaine nouveau qui se voit *aidé* par les récentes avancées dans l'apprentissage automatique [2].

Multimodal

Afin d'assurer une aisance d'utilisation, les SPAs sont fréquemment amenés à récupérer les requêtes (ou données) en entrée de la manière la plus naturelle possible (par exemple par le biais de la parole). Cependant, pour garantir une expérience d'utilisation adéquate, l'assistant sera souvent confronté à récupérer ces requêtes de différentes manières, que ce soit à travers une interface graphique (écran tactile) ou à travers un texte brut tapé au clavier, voire même à travers des expressions faciales ou des états cognitives/émotionnels [5], pour ensuite produire une réponse qui elle aussi pourrait éventuellement être de la forme textuelle, sonore ou les deux. Cette capacité à recevoir en entrée et/ou produire une sortie de plusieurs façons différentes est appelée la multi-modalité [6]. Cette caractéristique permet de masquer à l'utilisateur toute la complexité d'acquisition de ses requêtes.

Anthropomorphe

Plusieurs auteurs tendent à attribuer une grande importance à l'anthropomorphisme des SPAs [7], qui est

“Un mécanisme qui pousse les êtres humains à induire qu'une entité non-humanoïde possède des caractéristiques et comportements propres à l'homme”[8]

Ce comportement humanoïde pousserait donc l'utilisateur à se sentir plus à l'aise avec l'assistant, le conduisant ainsi à adopter une façon de communiquer plus humaine et moins structurée qu'avec les autres machines. Ceci est une caractéristique majeure d'un SPA se disant personnalisé.

Multi-plateforme et Flexible

Malgré leurs récentes prouesses, certains SPAs sont encore restreints à un écosystème fortement dépendant du fabricant. Cowan et al. mentionnent dans [9] que Apple Siri est limité à l'environnement constitué des produits de la firme à la pomme, n'ouvrant par défaut que les applications de cette dernière quand une requête lui est transmise. C'est un comportement que les assistants devraient éviter, car une indépendance des plateformes utilisées est, certes, très complexe à instaurer, mais offre plus de possibilités aux utilisateurs et aux développeurs pouvant ainsi exploiter la puissance de certaines plateformes (Smartphones, TV connectées, etc). Avec l'émergence de l'IoT (Internet of Things) et des maisons intelligentes par exemple, c'est un tout nouveau terrain de jeu qui est présenté aux SPAs, offrant plus d'opportunités pour les utilisateurs.

Domaines d'applications des SPAs

Après avoir vu les différents aspects que les SPAS doivent traiter, nous nous intéresserons maintenant aux types de services et applications que ces derniers pourraient fournir pour démontrer qu'ils peuvent bel et bien faciliter certaines tâches à l'homme.

Vie quotidienne

À la base, les SPAs étaient destinés à un usage très personnel comme la gestion des achats dans les supermarchés, ou des guides touristiques de plusieurs destinations de voyage. Cette spécificité a commencé à s'estomper petit à petit avec l'émergence de nouveaux systèmes dédiés à des applications plus générales, comme les maisons intelligentes ou les assistants de planification de tâches. Ceci a permis de mettre encore plus l'accent sur cet aspect de convivialité que les tout premiers SPAs ont tenté de perfectionner. Ainsi, ces assistants spécialisés dans des domaines restreints (Tourisme, shopping, détente, etc) ont été regroupés dans un seul système plus polyvalent, capable de répondre à des besoins quotidiens divers et variés, allant même à fournir une assistance aux personnes âgées pour leur faciliter les tâches rudimentaires devenues trop fatigantes.

Assistance professionnelle

Les SPAs ont aussi une place dans le monde professionnel. Dans [10] il est cité que dans les situations où la marge d'erreur est très petite (par exemple dans les système de manufacturing¹) l'assistance d'un SPA est nécessaire, servant d'un aide à l'humain pour la prise de décision.

Par exemple, dans un environnement de travail hétérogène (Nouveaux/anciens employés, Hiérarchies des postes ...) les SPAs pourraient décharger les employés les plus expérimentés de la tâche

1. Manufacturing ici dans le sens chaîne de montage industrielle, par exemple dans des usines.

d'assister les nouveaux arrivants, pour ainsi aider ces derniers dans leurs tâches et permettre aux autres de se focaliser sur les leurs.

E-Apprentissage

Les SPAs peuvent aussi être utiles dans l'enseignement, aussi bien dans un milieu académique que professionnel. D'une part ils pourraient occuper plusieurs rôles dans les établissements scolaires (correcteur automatique de copies, enseignant interactif ...) [11] et, d'autre part, accompagner les employés durant leurs formations professionnelles.

Ainsi, en considérant les caractéristiques d'un SPA, la sensibilité au contexte est reliée aux expériences antérieures de l'apprenant, permettant au SPA d'adapter son processus d'enseignement en conséquence.

En ce qui concerne l'aspect évolutif du SPA, il lui permet de préférer une approche d'enseignement à une autre selon les résultats de ses apprenants.

Exemples de SPAs

Pour illustrer la puissance des SPAs les plus récents, nous présentons dans cette section les quatre produits qui dominent le marché courant :

Google assistant

Lancé en 2016 sous forme d'un chatbot intégré dans l'application Google Allo, Google Assistant s'est vu ensuite être directement intégré sur les système d'exploitation Android (que ce soit sur smartphones ou tablettes, et plus récemment sur Google Home²). Google Assistant est un assistant à tout faire qui a été développé par les ingénieurs de Google dans le but de faciliter la recherche sur internet, la planification des tâches, l'ajustement des réglages de l'appareil, etc. Son point fort est sa capacité à engager une conversation bi-directionnelle avec l'utilisateur, assurant ainsi une interaction personnalisée variant d'un utilisateur à un autre. Cette capacité lui permet par exemple de proposer certains résultats de recherche selon les précédentes interactions avec l'utilisateur ou de lui proposer une activité si ce dernier lui mentionne qu'il s'ennuie (voir figure 2.1).

2. Appareil servant à contrôler les composants d'une smart-house ainsi que l'utilisation des différents services de Google

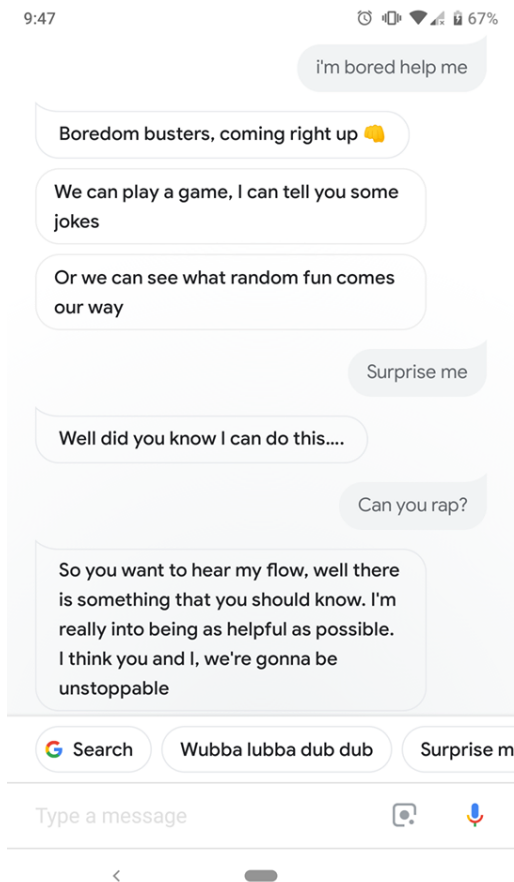


FIGURE 2.1 – *Conversation aléatoire avec Google Assistant*

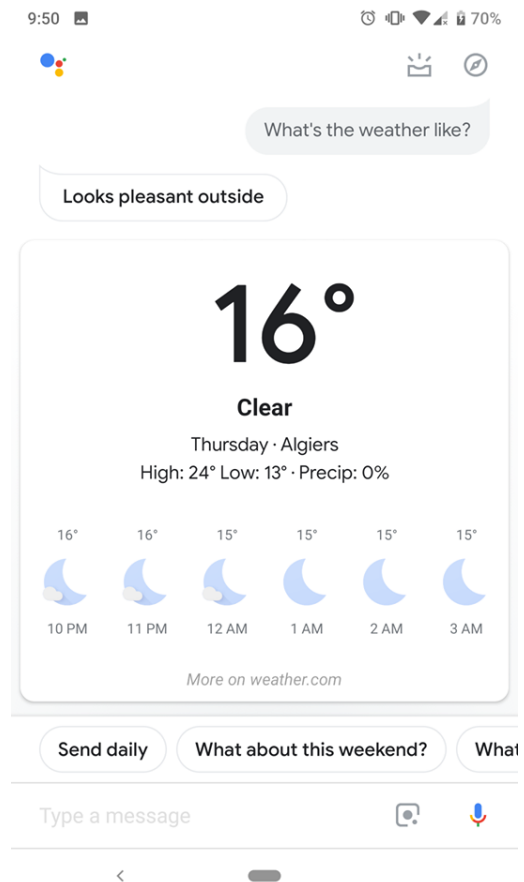


FIGURE 2.2 – *Requête simple formulée à Google Assitant*

Google duplex

Une des nouveautés impressionnante de Google Assistant est la fonctionnalité Google Duplex. Toujours en phase de développement, ce module est capable de passer des appels a de vraies personnes et d'avoir une conversation avec elles afin de réaliser une tâche demandée par l'utilisateur comme par exemple réserver une chambre d'hôtel, une table au restaurant, etc.

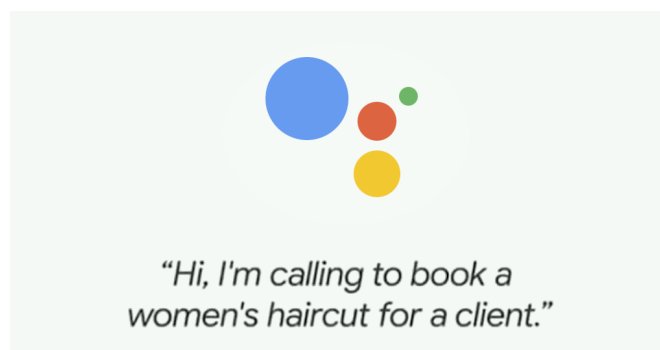


FIGURE 2.3 – *Google duplex réservant une place dans un salon de coiffure*

Apple Siri

Siri est l'assistant virtuel développé par Apple. Contrairement aux SPAs durant sa sortie, Siri proposait une nouvelle façon de communiquer avec l'utilisateur, à travers une interface de requêtes vocales, et une façon de **converser** très humanoïde (satisfaisant ainsi le critère d'anthropomorphisme 2.3.4). Siri est capable de répondre à des questions précises (voir figure 2.6), de proposer des recommandations, déléguer la requête à des services web ou d'autres applications (voir figures 2.4 et 2.5). Il a l'avantage (et l'inconvénient) d'être uniquement disponible que sur les appareils qui composent l'écosystème d'Apple (MacBook, iPhone, iWatch, etc).



FIGURE 2.4 – *Intégration aux applications [12]*

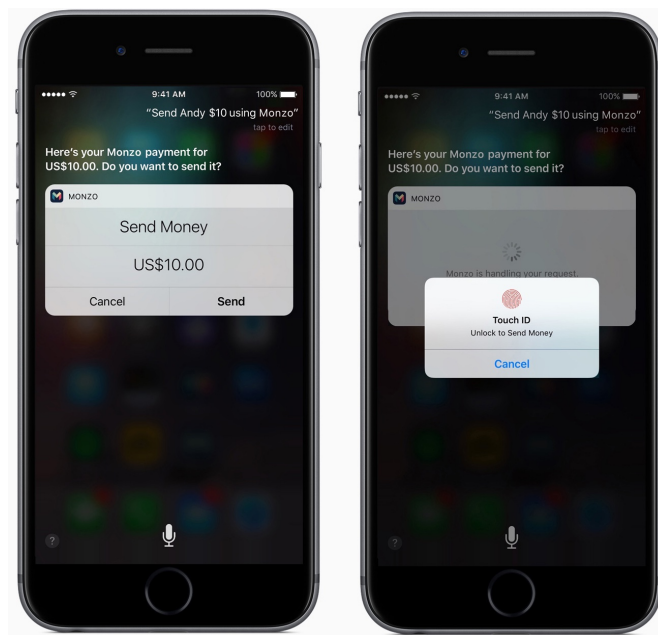


FIGURE 2.5 – *Service paiement 1 [12]*

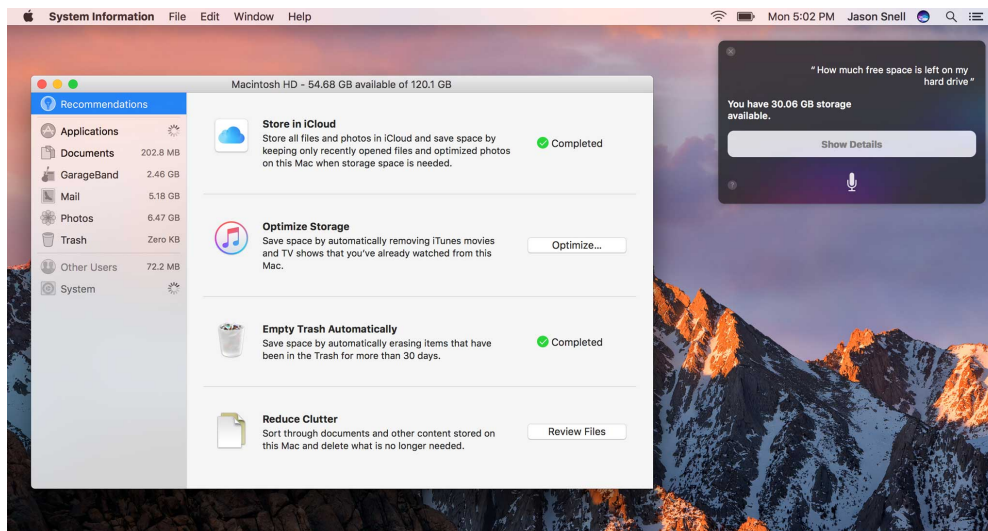


FIGURE 2.6 – Siri sur un laptop [13]

Amazon Alexa

Amazon Alexa est un assistant exclusivement intégré au dispositif Amazon Echo (un haut-parleur portatif). À l'instar de Siri, il est aussi capable de communiquer avec l'utilisateur par le biais de la parole, pouvant ainsi exécuter diverse commandes comme jouer de la musique, réciter des livres audios, annoncer des news en temps réel (Résultats sportifs, tendances politiques, etc). Son atout majeur est sa capacité à s'intégrer à plusieurs appareils-connectés (Contrôleur de thermostat ou de lumières ambiantes dans une Smart-House) ainsi que la possibilité d'ajout des Skills (ou compétences) de la part des développeurs tiers pour enrichir la panoplie de services que peut offrir Alexa.

Microsoft Cortana

Cortana est la tentative de la part de Microsoft d'intégrer un assistant dans son système d'exploitation Windows 10 et WindowsPhone. Il propose divers services de base tel que planifier des tâches, exécuter des commandes via la parole, et analyser des résultats de recherche sur le moteur de recherche de Microsoft, Bing, pour répondre à des questions.

Conclusion

À travers les sections précédentes, nous avons essayé de présenter les différents aspects d'un assistant virtuel intelligent (caractéristiques, exemples, architectures possibles, etc). Nous avons donc pu apprécier la potentielle puissance d'un tel système s'il venait à être perfectionner d'avantage.

En effet, en examinant les domaines d'applications, il est facile de déduire que le recours à un SPA peut grandement faciliter certaines tâches, que ce soit celles qui sont les plus triviales pouvant retarder d'autres tâches plus importantes, ou bien celles qui doivent faire appel à la précision

et à la grande capacité de calcul des machines, assurant ainsi des résultats précis et rapidement délivrés.

À la fin de ce chapitre nous pouvons donc mettre en valeur la place primordiale que pourraient avoir les SPAs s'ils arrivaient à maturité, c.à.d briser la barrière qui sépare les humains de la machine, parvenant ainsi à faire partie de la vie quotidienne des utilisateurs.

Dans le prochain chapitre nous allons principalement aborder les aspects techniques des différents composants du SPA que nous désirons réaliser.

Chapitre 3

Composants de base d'un assistant personnel

Introduction

Architecture d'un SPA

Reconnaissance automatique de la parole (ASR)

Définition

Modèle acoustique

Modèle de langage

Méthodes utilisées

N-grammes

Modèle de Markov Caché

Compréhension du langage naturel (NLU)

Définition

Classification d'Intent

Extraction d'entités

Analyse sémantique

Gestion du dialogue

La compréhension du langage naturel permet de transformer un texte en une représentation sémantique. Afin qu'un système puissent réaliser un dialogue aussi anthropomorphe que possible, il doit décider, à partir des représentations sémantiques reçues au cours du dialogue, quelle action à prendre à chaque étape de la conversation afin de la transmettre au générateur du langage naturel 3.6 pour afficher le résultat à l'utilisateur. On distingue deux principaux modules généralement présent dans les systèmes de gestion de dialogue :

- Un module qui met à jour l'état du gestionnaire de dialogue.
- Un module qui détermine la politique d'action du gestionnaire de dialogue.

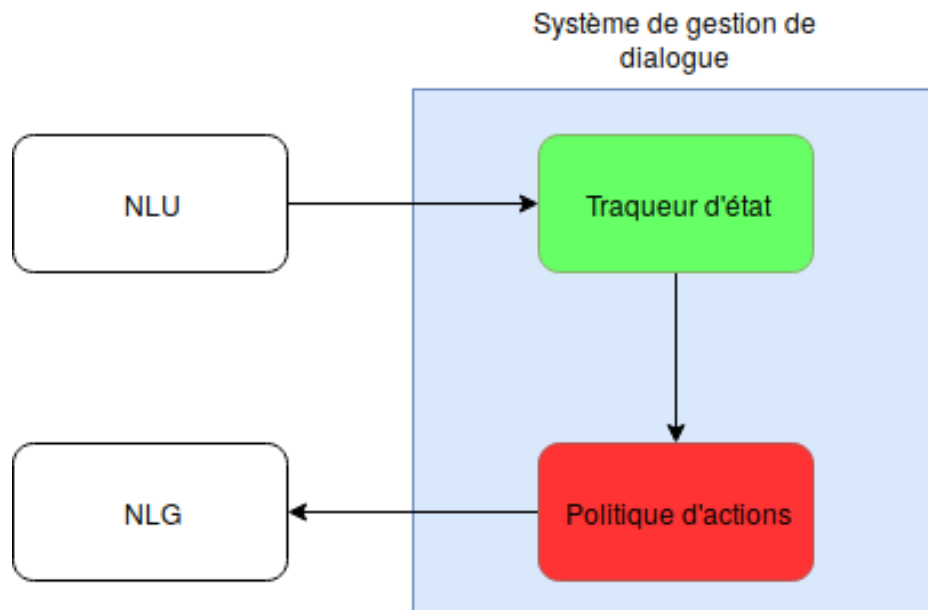


FIGURE 3.1 – Schéma général d'un gestionnaire de dialogue

Processus de décision Markovien (MDP)

Un gestionnaire de dialogue peut être modélisé par un processus de décision Markovien[14]. Ce dernier est modélisé par un 4-tuple $(S, A, P, R)^{(*)}$:

- S : ensemble d'états du système.
- A : ensemble d'actions du système.
- P : distribution de probabilités de transitions entre états sachant l'action prise. $P(s'/s, a)$ est la probabilité de passer à l'état s' sachant qu'on était à l'état s après avoir pris l'action a .
- R : est la récompense reçue immédiatement après avoir changer d'état avec une action donnée. $R(s'/s, a)$ est la récompense reçu après avoir passer à l'état s' sachant qu'on était à l'état s après avoir pris l'action a .

D'après $(*)$ un MDP, à tout instant t , est dans un état s , dans notre cas c'est l'état du gestionnaire de dialogue. Il peut prendre une action a afin de passer à un nouvel état s' , et sur ce fait, il reçoit une récompense, qui ,dans notre cas, est une mesure sur les performance du système de dialogue. Le but est de maximiser les récompense reçu.

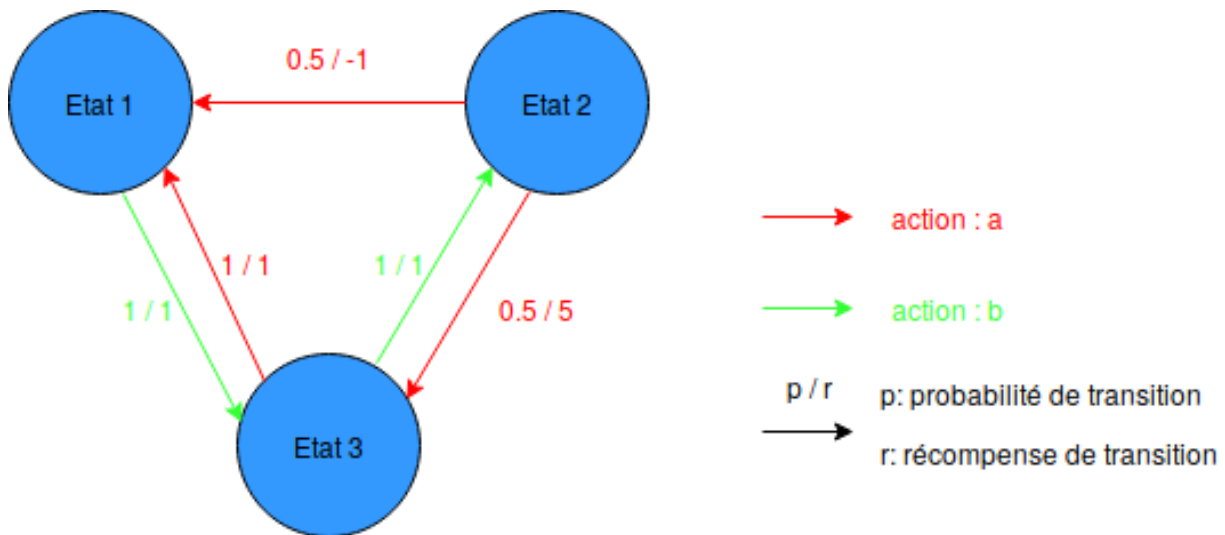


FIGURE 3.2 – Schéma représentant les transitions entre états dans un MDP

État du gestionnaire de dialogue

L'état d'un système de dialogue est une représentation sémantique qui contient des informations sur le but final de l'utilisateur ainsi que l'historique de la conversation. La représentation souvent utilisée dans les systèmes de dialogue est celle du cadre sémantique [15]. Cette structure contient des emplacements à remplir sur un domaine donné, la figure 3.3 illustre un exemple de cadre sémantique.

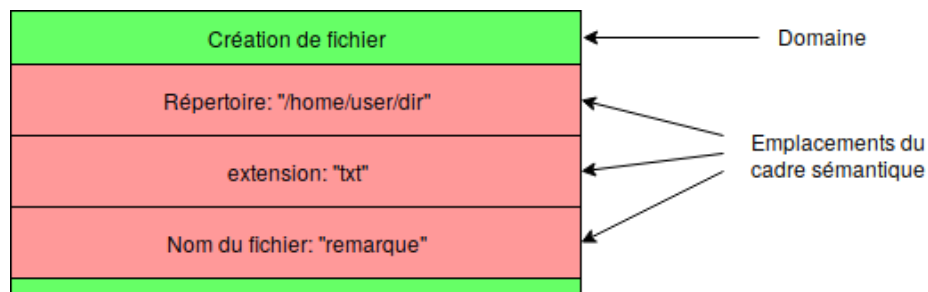


FIGURE 3.3 – Schéma représentant un cadre sémantique avec comme domaine : création de fichier

À l'arrivée d'une nouvelle information, un module dédié met à jour l'état du gestionnaire du dialogue. Comme l'action du système de dialogue est décidée à partir de son état, cette tâche est donc essentiel au bon fonctionnement du système. Plusieurs méthodes ont été donc proposé pour gérer le suivi de l'état du gestionnaire de dialogue.

Suivi de l'état avec une base de règles

La méthode traditionnelle utilisée est d'écrire manuellement les règles à suivre lors de l'arrivée d'une nouvelle information pour mettre à jour l'état[16]. Cependant, les bases de règles sont très susceptibles à faire des erreurs[15] comme ils sont moins robuste face aux incertitudes.

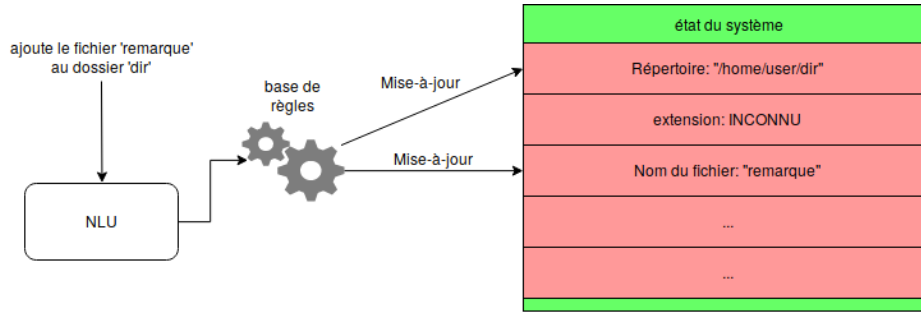


FIGURE 3.4 – Schéma représentant la mise-à-jour de l'état par un système basé règles

Suivi de l'état avec des méthodes statistiques

Le suivi dans ce cas se fait en gardant une distribution de probabilités sur l'état du système. D'où, l'utilisation des processus de décision markovien partiellement observé (POMDP)[17] qu'on introduira par la suite. Dans ce cas, le système garde une distribution de probabilités sur les valeurs possibles des différents emplacements du cadre sémantique.

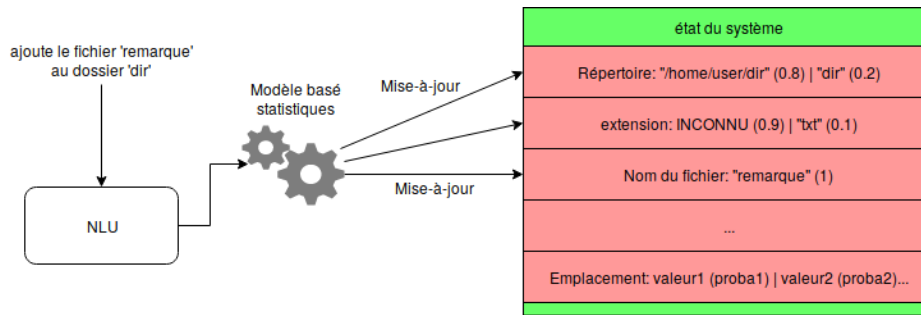


FIGURE 3.5 – Schéma représentant la mise-à-jour de l'état par un système basé statistiques

Processus de décision markovien partiellement observé (POMDP)

Comme dans les processus de décision markovien, un POMDP[18] passe d'un état à un autre en prenant une des actions possibles. Cependant, ce dernier ne connaît pas l'état exacte dans lequel il se trouve à un instant t . Il reçoit par contre une observation, dans notre cas c'est l'action de l'utilisateur, à partir de laquelle il peut estimer une distribution de probabilités sur l'état actuel. Pour résumer cela, un POMDP est un 6-tuple (S, A, P, R, M, O) :

- Les 4 premiers composants sont les même que celui d'un MDP 3.5.1.
- M : l'ensemble des observation.
- O : distribution de probabilités sur les observations o sachant en connaissant l'état et l'action prise pour y arriver. $O(o|s,a)$ est la probabilité d'observer o sachant qu'on se trouve à l'état s et qu'on a pris l'action a pour y arriver.

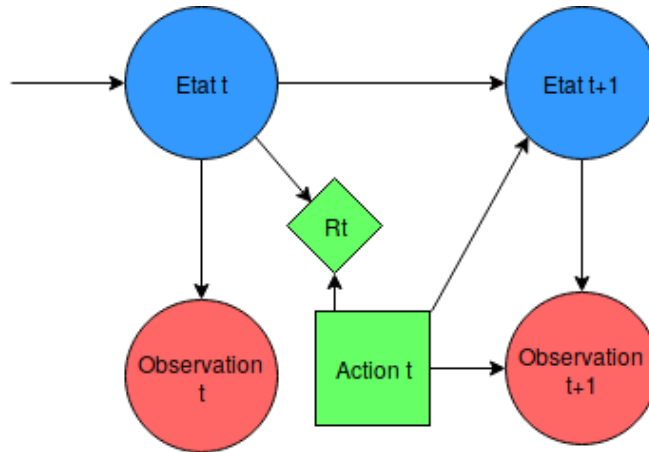


FIGURE 3.6 – *Diagramme d'influence dans un POMDP*

Suivi de l'état avec réseaux de neurones profonds

Récemment, des approches utilisant les réseaux de neurones profonds (refpart2) ont fait leurs apparitions. En effet, l'utilisation des architectures profondes permet de capter des relations complexes entre les caractéristiques d'un dialogue et ainsi mieux estimer l'état du système. Le réseau de neurones estime les probabilités de toutes les valeurs possibles d'un emplacement du cadre sémantique[19]. En conséquence, il peut être utilisé comme modèle de suivi d'état pour un processus partiellement observable.

Politique de gestion de dialogue

La première partie était dédiée au module qui suit l'état du système de dialogue. Dans cette partie, Nous allons présenter des approches proposées afin d'arriver au but du MDP, c'est à dire quelles actions prendre pour maximiser la somme des récompenses obtenus.

Gestion de dialogue avec une base de règles

Les premières approches utilisaient des systèmes de règles destiné à un domaine bien spécifique. Elles étaient déployés dans plusieurs domaines d'application pour sa simplicité. Cependant, le travail manuel nécessaire reste difficile à faire, et, généralement, n'aboutit pas à des résultats flexibles qui peuvent suivre le flux du dialogue convenablement[20].

Gestion de dialogue par apprentissage

La résolution d'un MDP revient à trouver une estimation de la fonction de récompense afin de pouvoir choisir la meilleure action. La majorité des approches récentes utilise l'apprentissage par renforcement pour but d'estimer la récompense obtenue par une action et un état donnés. Cette préférence par rapport aux approches supervisées revient à la difficulté de produire des corpus de dialogues[21], encore moins des corpus annotés avec les récompenses à chaque transition. Néanmoins, il existe des approches de bout en bout qui exploite des architectures avec réseaux de

neurones profonds et traite le problème comme Seq2Seq(part2) afin de produire directement une sortie à partir des informations reçues par l'utilisateur[22][40].

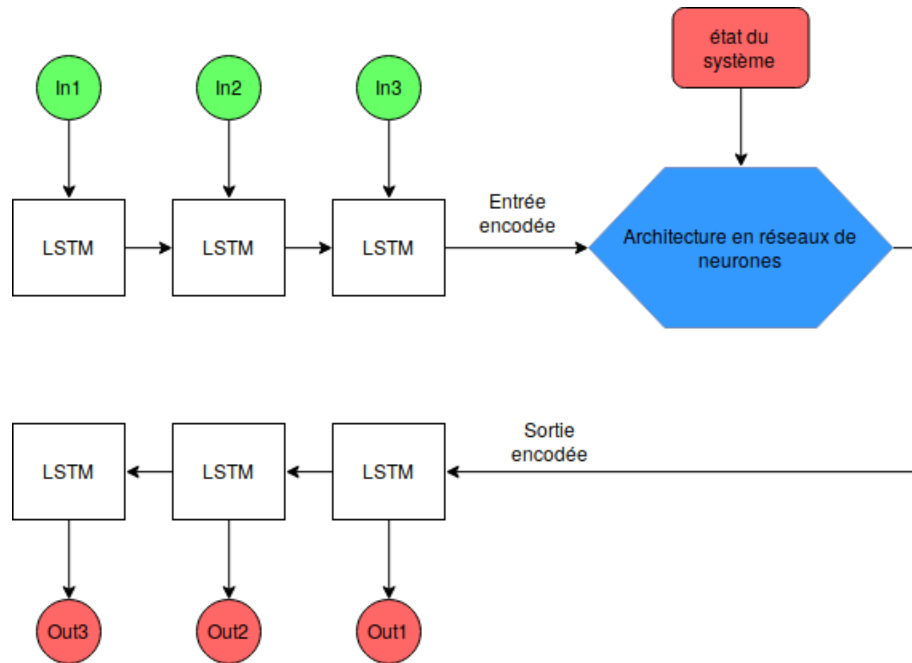


FIGURE 3.7 – Schéma de gestion de dialogue de bout en bout avec architecture Seq2Seq

Apprentissage par renforcement

Simulateur d'utilisateur

Génération du langage naturel (NLG)

Le domaine de la génération automatique du langage naturel est l'un des domaines dont les bordures sont difficiles à définir (Evans et al., 2002). il est vrai que la sortie d'un tel système est clairement du texte. Cependant, l'ambiguïté se trouve dans ses l'entrées, c'est à dire, sur quoi se basera le système pour générer le texte. D'après (Reiter & Dale, 1997)[23] la génération du langage naturel est décrite comme étant le sous domaine de l'intelligence artificielle qui traite la construction des systèmes de génération de texte à partir d'une représentation non-linguistique de l'information, celle-ci peut être une représentation sémantique, des données numériques, une base de connaissances ou même des données visuelles (images ou vidéos). Ceci dit, d'autres travaux ,comme Labbé & Portet (2012)[24], utilisent les même techniques pour des entrées linguistique. Enfin la génération du langage naturel peut être très proche de la gestion de dialogue[25], en effet, le texte généré doit prendre en compte l'historique de la conversation et le contexte de l'utilisateur. Il existe six tâches trouvées fréquemment dans les systèmes de génération de texte [23].

Détermination du contenu

Cette partie consiste à sélectionner les informations de l'entrée dont le système veut transmettre le contenu sous forme de texte naturel à l'utilisateur. En effet, les données en entrée peuvent contenir plus d'informations que ce que l'on désire communiquer[26], de plus, cette information peut aussi dépendre de l'utilisateur et de ses connaissances[25]. Ce qui requiert de mettre au point un système qui détecte les informations pertinentes à l'utilisateur.

Structuration de texte

Après la détermination du contenu, le système doit ordonner les informations à transmettre. Ceci dépend grandement du domaine d'application qui peut exiger des contraintes d'ordre temporelle ou de préférence par importance des idées. Les informations à transmettre en elles-mêmes sont souvent reliées par sens ce qui implique une certaine structuration de texte à respecter.

Agrégation de phrases

Certaines informations peuvent être transmises dans une même phrase. Cette partie introduit des notions de la linguistique afin que le texte généré soit plus lisible et éviter les répétition. Un exemple de cela peut être la description de la météo à Alger au cours de la matinée :

- Il va faire 16° à Alger à 7h.
- Il va faire 17° à Alger à 8h.
- Il va faire 18° à Alger à 9h.
- Il va faire 18° à Alger à 10h.

Ceci peut être agrégé en un texte plus compacte : "La température moyenne à Alger sera de 17° entre 7h et 10h."

Lexicalisation

Le système choisit les mots et les expressions à utiliser pour communiquer le contenu des phrases sélectionnées. La difficulté de cette tâche revient à l'existence de plusieurs manières d'exprimer la même idée. Cependant, certains mots ou expressions sont plus appropriés en certaines situations que d'autres. En effet, "inscrire un but" est une façon inadéquate d'exprimer un but contre son camp[27].

Génération d'expressions référentielles (REG)

Cette partie du système se focalise dans la génération d'expressions référentielles qui peuvent être entre autres : noms propres, groupes nominaux ou pronoms et ceci a pour but d'identifier les entités du domaine. Cette tâche semble être très proche de sa prédécesseur ; elle s'avère néanmoins plus délicate dû à la difficulté de confier suffisamment d'information sur l'entité afin de la différencier des autres[23]. Le système doit faire un choix de l'expression référentielle en se basant sur plusieurs facteurs, par exemple "Mohammed", "Le professeur" ou "Il" font référence à la même personne. Cependant, le choix entre eux dépendrait de si l'entité a été mentionnée auparavant et des détails l'accompagnant par exemple.

Réalisation linguistique

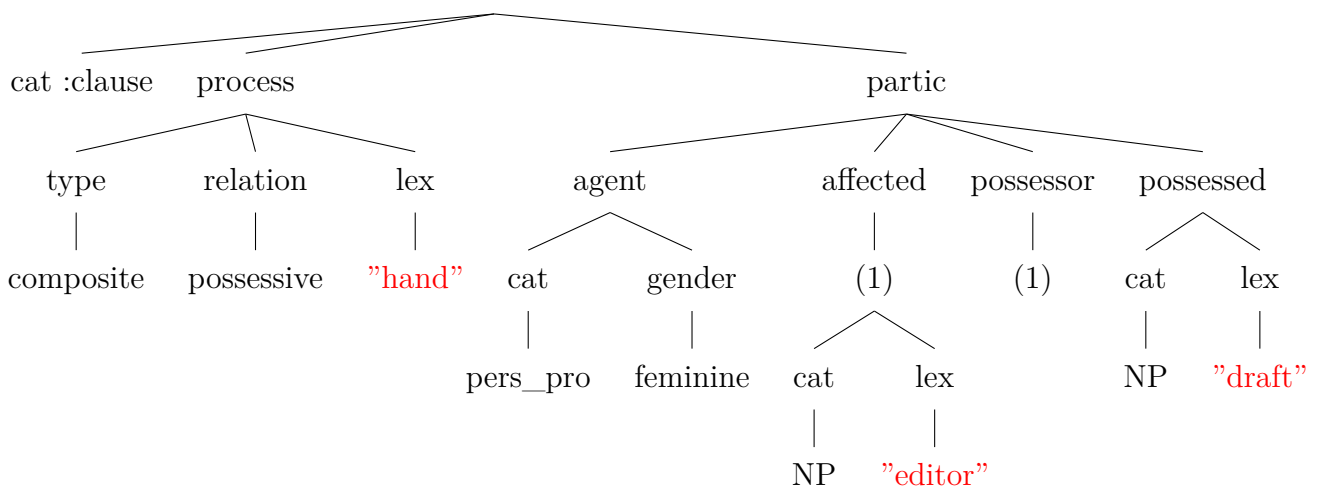
La dernière tâche consiste à combiner les mots et expressions sélectionnés pour construire une phrase linguistiquement correcte. Ceci requiert l'utilisation des bonnes formes morphologiques des mots, les ordonner et éventuellement l'addition de certains mots du langage afin de réaliser une structure de phrase grammaticalement et sémantiquement correcte. Plusieurs méthodes ont été proposées, principalement les méthodes basées sur des règles manuellement construites (modèles de phrases, systèmes basés grammaires) ou des approches statistiques [27].

Modèles de phrases : La réalisation se fait en utilisant des modèles de phrases prédéfinis. Il suffit de remplacer des espaces réservés par certaines entrées du système. Par exemple, une application dans un contexte météorologique pourrait utiliser le modèle suivant : "la température à

[ville] atteint [température]° le [date]”.

Cette méthode est utilisée lorsque les variations des sorties de l’application sont minimales. Son utilisation a l’avantage et l’inconvénient d’être rigide. D’un coté il est facile de contrôler la qualité des sorties syntaxiquement et sémantiquement tout en utilisant des règles de remplissage complexe [28]. Cependant, lorsque le domaine d’application présente beaucoup d’incertitude, cette méthode exige un travail manuel énorme, voire impossible à faire, pour réaliser une tâche pareille. Bien que certains travaux ont essayer de faire un apprentissage de modèles de phrases à partir d’un corpus[29] cette méthode reste inefficace lorsqu’il s’agit d’application qui nécessite un grand nombre de variations linguistiques.

Systèmes basés grammaire : La réalisation peut se faire en suivant une grammaire du langage. Celle-ci contient les règles morphologiques et de structures de la langues, notamment la grammaire systémique fonctionnelle (SFG)[30] a été largement utilisé comme dans NIGEL[31] ou KPML[32]. L’exploitation des grammaires dans la génération du texte nécessite généralement des entrées détaillées. En plus des composantes du lexique sélectionnées, des descriptions de leurs rôles ainsi que leurs fonctions grammaticales sont souvent exigées. Un exemple d’entrée d’un système basé grammaire est celui de SURGE[33] :



Qui génère la phrase : “She hands the draft to the editor”.

Comme les modèles de phrases, les systèmes basés grammaire nécessite un énorme travail manuel. En particulier, il est difficile de prendre en compte le contexte en définissant les règles de choix entre les variantes possibles du texte résultat à partir des entrées[27].

Approches statistiques : Il existe plusieurs méthodes basées sur des statistiques pour la tâche de réalisation. Certains se basent sur des grammaires probabilistes, cette dernière a l’avantage de minimiser le travail manuel tout en couvrant plus de cas de réalisation. Il existe principalement deux approches l’utilisant[27] :

- La première se base sur une petite grammaire qui génère plusieurs alternatives qui sont ensuite ordonnés selon un modèle statistique basé sur un corpus pour sélectionner la phrase la plus probable (par exemple Langkilde-Geary (2000)[34]).

- La deuxième méthode utilise les informations statistiques directement au niveau de la génération pour produire la solution optimale (exemple : Belz (2008)[35]).

Dans les deux méthodes sus-citées la grammaire de base peut être manuellement faite, dans ce cas, les informations statistiques aideront à la détermination de la solution optimale, ou elle peut être extraite à partir des données, comme l'utilisation des Treebanks¹ pour déduire les règles de grammaire[36].

D'autres approches statistiques n'utilisent pas des grammaires mais se basent sur des classificateurs. Ces derniers peuvent être cascades de telle sorte à décider quel constituant utiliser dans quelle position ainsi que les modifications nécessaires pour générer un texte correcte. À noter qu'une telle approche, ne nécessitant pas l'utilisation de grammaire, utilise des entrées plus abstraites et moins détaillées linguistiquement. À voir même la possibilité de s'étendre aux autres tâches de NLG, c'est à dire un système qui accomplit plusieurs tâches de NLG en parallèle en utilisant les entrées initiales. Dans la suite de ce travail nous allons présenter certains de ces systèmes qui sont plus utilisés récemment.

Systèmes basés encodeur-décodeur

Une architecture souvent utilisée dans le traitement du langage naturel est l'encodeur-décodeur (part2). En particulier, son utilisation dans les tâches seq2seq (part2) ce qui permet de mettre en correspondance une séquence de taille variable en entrée avec une autre séquence en sortie. Les modèles seq2seq peuvent être adapter pour convertir une représentation abstraite de l'information en langage naturel[37].

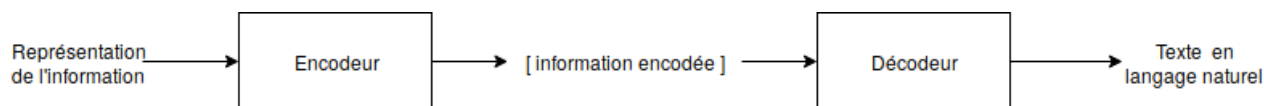


FIGURE 3.8 – Schéma d'une architecture encodeur-décodeur pour NLG

Beaucoup d'approches de génération de langage naturel en gestion de dialogue utilise des encodeur-décodeur. Wen et al. (2015)[38] utilise par exemple des LSTMs(part2) sémantiquement conditionnés ; il ajoute aux LSTMs classiques une couche contenant des informations sur l'action prise par le gestionnaire de dialogue pour assurer que la génération représente le sens désiré. D'autres travaux utilisent des réseaux de neurones récurrent (part2) pour encoder l'état du gestionnaire de dialogue et l'entrée reçu suivis par un décodeur pour générer le texte de la réponse[39][40][41].

1. un Treebank est un texte analysé qui contient des informations syntaxiques ou sémantiques sur les structures de phrases

Table des figures

2.1	Conversation aléatoire avec Google Assistant	10
2.2	Requête simple formulée à Google Assitant	10
2.3	Google duplex réservant une place dans un salon de coiffure	10
2.4	Intégration aux applications [12]	11
2.5	Service paiement 1 [12]	11
2.6	Siri sur un laptop [13]	12
3.1	Schéma général d'un gestionnaire de dialogue	16
3.2	Schéma représentant les transitions entre états dans un MDP	17
3.3	Schéma représentant un cadre sémantique avec comme domaine : création de fichier	17
3.4	Schéma représentant la mise-à-jour de l'état par un système basé règles	18
3.5	Schéma représentant la mise-à-jour de l'état par un système basé statistiques	18
3.6	Diagramme d'influence dans un POMDP	19
3.7	Schéma de gestion de dialogue de bout en bout avec architecture seq2seq	20
3.8	Schéma d'une architecture encodeur-décodeur pour NLG	24

Bibliographie

- [1] D. A. Norman, *The design of everyday things*. New York : Basic Books, 2002.
- [2] R. Knote, A. Janson, L. Eigenbrod, and M. Söllner, “The what and how of smart personal assistants : Principles and application domains for is research,” in *Multikonferenz Wirtschaftsinformatik (MKWI)*, 2018.
- [3] H. Gellersen, *Handheld and Ubiquitous Computing : First International Symposium, HUC’99, Karlsruhe, Germany, September 27-29, 1999, Proceedings (Lecture Notes in Computer Science)*. Springer, 1999.
- [4] S. J. Russell and P. Norvig, *Artificial Intelligence : A Modern Approach*. Pearson Education, 2 ed., 2003.
- [5] T. Dingler, “Cognition-aware systems as mobile personal assistants,” in *Proceedings of the 2016 ACM International Joint Conference on Pervasive and Ubiquitous Computing Adjunct - UbiComp ’16*, ACM Press, 2016.
- [6] E. Luger and A. Sellen, “”like having a really bad PA”,” in *Proceedings of the 2016 CHI Conference on Human Factors in Computing Systems - CHI ’16*, ACM Press, 2016.
- [7] R. Trappl, ed., *Your Virtual Butler*. Springer Berlin Heidelberg, 2013.
- [8] A. Purington, J. G. Taft, S. Sannon, N. N. Bazarova, and S. H. Taylor, “”alexa is my new BFF”,” in *Proceedings of the 2017 CHI Conference Extended Abstracts on Human Factors in Computing Systems - CHI EA ’17*, ACM Press, 2017.
- [9] B. R. Cowan, N. Pantidi, D. Coyle, K. Morrissey, P. Clarke, S. Al-Shehri, D. Earley, and N. Bandeira, “”what can i help you with?”,” in *Proceedings of the 19th International Conference on Human-Computer Interaction with Mobile Devices and Services - MobileHCI ’17*, ACM Press, 2017.
- [10] J. Imtiaz, N. Koch, H. Flatt, J. Jasperneite, M. Voit, and F. van de Camp, “A flexible context-aware assistance system for industrial applications using camera based localization,” in *Proceedings of the 2014 IEEE Emerging Technology and Factory Automation (ETFA)*, IEEE, sep 2014.
- [11] “Engaging the appropriation of technology-mediated learning services – a theory-driven design approach research in progress,” 2015.
- [12] “Apple shares examples of siri’s third-party app integration on ios 10.” <https://www.idownloadblog.com/2016/09/01/apple-siri-ios-10-app-integration/>. (Accessed on 10/29/2018).
- [13] “macos sierra review : Hey siri, where did my files go? - six colors.” <https://sixcolors.com/post/2016/09/sierra-review/>, 2016. (Accessed on 10/29/2018).
- [14] t. . A. M. D. P. Richard Bellman” *Indiana Univ. Math. J.*, *fjournal* = ”Indiana University Mathematics Journal”, vol. 6, pp. 679–684, 1957.

- [15] H. Chen, X. Liu, D. Yin, and J. Tang, “A survey on dialogue systems : Recent advances and new frontiers,” *SIGKDD Explor. Newsl.*, vol. 19, pp. 25–35, Nov. 2017.
- [16] D. Goddeau, H. Meng, J. Polifroni, S. Seneff, and S. Busayapongchai vol. 2, 11 1996.
- [17] S. Young, M. Gašić, S. Keizer, F. Mairesse, J. Schatzmann, B. Thomson, and K. Yu, “The hidden information state model : A practical framework for pomdp-based spoken dialogue management,” *Comput. Speech Lang.*, vol. 24, pp. 150–174, Apr. 2010.
- [18] K. J. Åström, “Optimal control of Markov Processes with incomplete state information,” *Journal of Mathematical Analysis and Applications*, vol. 10, pp. 174–205, January 1965.
- [19] M. Henderson, B. Thomson, and S. J. Young, “Deep neural network approach for the dialog state tracking challenge,” in *SIGDIAL Conference*, pp. 467–471, 2013.
- [20] C. Lee, S. Jung, K. Kim, D. Lee, and G. G. Lee, “Recent approaches to dialog management for spoken dialog systems,” *JCSE*, vol. 4, pp. 1–22, 2010.
- [21] J. Henderson, O. Lemon, and K. Georgila, “Hybrid reinforcement/supervised learning of dialogue policies from fixed data sets,” *Comput. Linguist.*, vol. 34, pp. 487–511, Dec. 2008.
- [22] T.-H. Wen, M. Gasic, N. Mrksic, L. M. Rojas-Barahona, P. hao Su, S. Ultes, D. Vandyke, and S. J. Young, “A network-based end-to-end trainable task-oriented dialogue system,” in *EACL*, pp. 438–449, 2017.
- [23] E. Reiter and R. Dale, “Building applied natural language generation systems,” *Natural Language Engineering*, vol. 3, pp. 57–87, Mar. 1997.
- [24] C. Labbé and F. Portet, “Towards an abstractive opinion summarisation of multiple reviews in the tourism domain,” *CEUR Workshop Proceedings*, vol. 917, pp. 87–94, 01 2012.
- [25] N. Dethlefs, “Context-sensitive natural language generation : From knowledge-driven to data-driven techniques,” *Language and Linguistics Compass*, vol. 8, pp. 99–115, 03 2014.
- [26] J. Yu, E. Reiter, J. Hunter, and C. Mellish, “Choosing the content of textual summaries of large time-series data sets,” *Natural Language Engineering*, vol. 13, pp. 25–49, Mar. 2007.
- [27] A. Gatt and E. Krahmer, “Survey of the state of the art in natural language generation : Core tasks, applications and evaluation,” *J. Artif. Int. Res.*, vol. 61, pp. 65–170, Jan. 2018.
- [28] M. Theune, E. Klabbers, J. R. De Pijper, E. Krahmer, and J. Odijk, “From data to speech : A general approach,” *Natural Language Engineering*, vol. 7, pp. 47–86, Mar. 2001.
- [29] G. Angeli, C. D. Manning, and D. Jurafsky, “Parsing time : Learning to interpret time expressions,” in *Proceedings of the 2012 Conference of the North American Chapter of the Association for Computational Linguistics : Human Language Technologies, NAACL HLT ’12*, (Stroudsburg, PA, USA), pp. 446–455, Association for Computational Linguistics, 2012.
- [30] M. A. K. Halliday and C. M. I. M. Matthiessen, *Introduction to Functional Grammar*. London : Hodder Arnold, 3 ed., 2004.
- [31] W. C. Mann and C. M. I. M. Matthiessen, “Nigel : A systemic grammar for text generation.,” 1983.
- [32] J. A. Bateman, “Enabling technology for multilingual natural language generation : The kpml development environment,” *Nat. Lang. Eng.*, vol. 3, pp. 15–55, Mar. 1997.
- [33] M. Elhadad and J. Robin, “An overview of surge : a reusable comprehensive syntactic realization component,” in *International Natural Language Generation Workshop*, pp. 1–4, 1996.
- [34] I. Langkilde-Geary, “Forest-based statistical sentence generation,” in *ANLP*, pp. 170–177, 2000.

- [35] A. Belz, “Automatic generation of weather forecast texts using comprehensive probabilistic generation-space models,” *Nat. Lang. Eng.*, vol. 14, pp. 431–455, Oct. 2008.
- [36] D. Espinosa, M. White, and D. Mehay, “Hypertagging : Supertagging for surface realization with ccg,” in *ACL*, pp. 183–191, 2008.
- [37] T. C. Ferreira, I. Calixto, S. Wubben, and E. Krahmer, “Linguistic realisation as machine translation : Comparing different mt models for amr-to-text generation,” in *INLG*, pp. 1–10, 2017.
- [38] T.-H. Wen, M. Gasic, N. Mrksic, P. hao Su, D. Vandyke, and S. J. Young, “Semantically conditioned lstm-based natural language generation for spoken dialogue systems,” in *EMNLP*, pp. 1711–1721, 2015.
- [39] A. Sordoni, M. Galley, M. Auli, C. Brockett, Y. Ji, M. Mitchell, J.-Y. Nie, J. Gao, and W. B. Dolan, “A neural network approach to context-sensitive generation of conversational responses,” in *HLT-NAACL*, pp. 196–205, 2015.
- [40] I. V. Serban, A. Sordoni, Y. Bengio, A. Courville, and J. Pineau, “Building end-to-end dialogue systems using generative hierarchical neural network models,” in *Proceedings of the Thirtieth AAAI Conference on Artificial Intelligence*, AAAI, pp. 3776–3783, AAAI Press, 2016.
- [41] R. Goyal, M. Dymetman, and E. Gaussier, “Natural language generation through character-based rnns with finite-state prior knowledge,” in *COLING*, pp. 1083–1092, 2016.