

Article

« Présentation d'un modèle de génération automatique »

Laurence Danlos

Revue québécoise de linguistique, vol. 13, n° 1, 1983, p. 203-228.

Pour citer cet article, utiliser l'information suivante :

URI: <http://id.erudit.org/iderudit/602510ar>

DOI: 10.7202/602510ar

Note : les règles d'écriture des références bibliographiques peuvent varier selon les différents domaines du savoir.

Ce document est protégé par la loi sur le droit d'auteur. L'utilisation des services d'Érudit (y compris la reproduction) est assujettie à sa politique d'utilisation que vous pouvez consulter à l'URI <https://apropos.erudit.org/fr/usagers/politique-dutilisation/>

Érudit est un consortium interuniversitaire sans but lucratif composé de l'Université de Montréal, l'Université Laval et l'Université du Québec à Montréal. Il a pour mission la promotion et la valorisation de la recherche. Érudit offre des services d'édition numérique de documents scientifiques depuis 1998.

Pour communiquer avec les responsables d'Érudit : info@erudit.org

PRÉSENTATION D'UN MODÈLE DE GÉNÉRATION AUTOMATIQUE

Laurence Danlos

Dans la perspective d'utiliser le langage naturel comme moyen de communication entre l'homme et la machine, notre travail porte sur la production de textes en langues naturelles. Le modèle que nous avons développé a pour but que la machine communique à l'utilisateur les informations qu'il désire dans la langue de son choix.

1. Présentation du modèle de génération

La tâche effectuée par notre système consiste à "traduire" en langue naturelle un ensemble de données représentées dans un formalisme abstrait. Ce formalisme abstrait est indépendant de la langue: aucune information linguistique n'y figure. De ce fait, la banque de données peut servir à d'autres fins que la génération de textes et elle peut être consultée par des utilisateurs parlant différentes langues.

La production d'un texte repose au moins sur deux types de décisions: le premier type de décisions, que nous appellerons "conceptuelles", concerne des problèmes comme: dans quel ordre doivent apparaître les informations? Quelles sont les informations qui doivent être exprimées explicitement par rapport à celles qui ne sont exprimées qu'implicitement? Le deuxième type de décisions, que nous appellerons "linguistiques", concerne des problèmes comme: comment découper le texte en paragraphes et en phrases? Quelles constructions syntaxiques choisir? Quels mots choisir?

La première question que l'on doit se poser est la suivante: est-ce que ces deux types de décisions peuvent être prises indépendamment? L'originalité de notre travail par rapport à d'autres recherches dans ce domaine est que nous répondons par la négative à cette question: ces deux types de décisions doivent être prises simultanément. À titre de brève illustration, considérons la relation causale suivante dont la cause est un acte et le résultat un état:

ACTE : TIR
 AGENT =: Max
 OBJET =: Luc

(A) CAUSE

ÉTAT : MORT
 OBJET =: Luc

Supposons que les décisions conceptuelles conduisent à exprimer d'abord la mort de Luc, ensuite l'acte de Max. De telles décisions, transmises à un module linguistique français, aboutiraient à des formes comme :

Luc est mort parce que Max lui a tiré dessus.

Luc est mort. Max lui a tiré dessus.

Or, dans de nombreux contextes, on préférera utiliser des formes telles que :

(1) a. Max a tué Luc en lui tirant dessus.

b. Luc a été tué par Max qui lui a tiré dessus.

construites autour du verbe tuer. Ces formes n'obéissent pas aux décisions conceptuelles dissociant la mort de Luc et l'acte de Max: *tuer* (dans la construction $N_0 V N_1$ =: *Max a tué Luc*¹) exprime en même temps la mort de N_1 et le fait que cette mort soit due à un acte (non précisé) de N_0 (McCawley, 1971)². Ces faits indiquent donc que les décisions conceptuelles ne doivent pas être prises sans considérer les possibilités offertes par la langue. Dans le cas présent, il faut prendre en compte l'existence de verbes à sémantique causative tels que *tuer*. De plus, s'il existait un niveau conceptuel indépendant

1. Dans cette notation, N_0 désigne le sujet, N_1 le premier complément et N_2 le deuxième.

2. Nous ne tenons pas compte ici de l'emploi "psychologique" de *tuer*: *Ces enfants me tuent*.

de toute considération linguistique, ce niveau conceptuel devrait permettre de fournir des résultats satisfaisants dans plusieurs langues. Or, en anglais, les formulations les plus appropriées pour (A) sont des formes telles que:

- (2) a. Max shot and killed Luc.
 b. Max shot Luc dead.

qui n'ont pas d'équivalent littéral en français. Il ne peut donc pas exister de décisions conceptuelles qui permettent de produire (1) en français et (2) en anglais. Ajoutons que les formes (2) sont quasiment figées: elles mettent en jeu des constructions syntaxiques qui sont distributionnellement très contraintes (Green, 1972):

- ?*Max shot and paralysed Luc.
 ?*Max shot Luc paralysed.

Elles ne peuvent donc pas être obtenues par des règles générales associant des combinaisons d'éléments de sens à des constructions syntaxiques.

Considérons maintenant le cas d'une relation causale de type (A) où l'ACTE est une strangulation au lieu d'un coup de feu. Cette relation causale est exprimée d'une façon satisfaisante par la forme:

Max a étranglé Luc.

où la mort de Luc est implicite: le verbe *étrangler*, dans la

construction N_0 *étrangler* N_1 avec un sujet N_0 humain³, désigne un acte de N_0 et implique la mort de N_1 sauf mention implicite du contraire. Aucune décision conceptuelle raisonnable ne peut aboutir au résultat de ne pas exprimer la mort de Luc. Il est donc sans intérêt de postuler l'existence d'un niveau conceptuel indépendant de la langue, puisque les résultats d'un tel module peuvent être contrecarrés par des faits linguistiques, tels que l'existence du verbe *étrangler*.

Enfin, comparons les paires suivantes:

- (1) a. Luc s'est tué en s'ouvrant les veines.
 b. Luc s'est ouvert les veines, se tuant.
- (2) a. Luc s'est suicidé en s'ouvrant les veines.
 b. *Luc s'est ouvert les veines, se suicidant.

L'inacceptabilité de (2b) montre que le choix entre la construction (a) (où la cause suit le résultat) et la construction (b) (où la cause précède le résultat) et le choix entre *se tuer* et *se suicider* sont des décisions qui doivent être prises simultanément. Autrement dit, les décisions sémantiques ne doivent pas être prises sans tenir compte des contraintes présentées par les mots, telles que l'interdiction pour *se suicider* d'apparaître dans la construction (b).

3. Le verbe *étrangler* se construit aussi avec un sujet non humain: *Cette chemise m'étrangle*. Il n'implique pas alors la mort de l'objet.

Ces quelques exemples illustrent un fait du langage naturel, à savoir qu'il n'existe pas de correspondance biunivoque entre éléments de sens et formes de phrases. Par cela, nous entendons qu'il n'est généralement pas possible de caractériser une notion sémantique par une ou plusieurs propriétés syntaxiques ni d'associer à une propriété syntaxique un ou plusieurs éléments de sens⁴. Cette caractéristique du langage naturel a été mise en évidence dans la grammaire développée par Z.S. Harris pour l'anglais et dans celle développée par M. Gross et ses collègues pour les langues romanes⁵. Ces grammaires, construites avec des méthodes propres à la linguistique, couvrent l'ensemble des phénomènes des langues étudiées; elles ne se bornent pas à des sous-ensembles de la langue, comme c'est souvent le cas des grammaires empruntant leurs méthodes à la logique ou à l'informatique, pour lesquelles les correspondances établies entre le sens et la forme pourraient n'être qu'un artefact dû au choix limité des exemples traités.

-
4. Cette affirmation dépasse donc la simple constatation qu'il n'existe pas de relation bijective entre la forme et le sens dans la mesure où à une forme peut correspondre plusieurs éléments de sens (phénomènes des homonymes et des ambiguïtés syntaxiques) et à un élément de sens peut correspondre plusieurs formes (phénomènes des synonymes et des paraphrases).
5. Citons entre autres Harris (1982), et les études de Gross (1981), Danlos (1981), Giry-Schneider (1981) et Guillet et Leclère (1981) sur la relation entre formes syntaxiques et prédicats sémantiques.

En plus des décisions conceptuelles et linguistiques évoquées, la production de textes demande que soient appliquées des règles de syntaxe telles que l'accord entre le sujet et le verbe ou la réduction d'une complétive à une infinitive:

Max aime que Luc chante.

*Max aime qu'il chante⁶.

=Max aime chanter.

Il est bien évident que ces règles de syntaxe sont indépendantes du type de données concernées et que ce point doit être reflété dans un système de génération.

Notre générateur est donc composé de deux modules: le premier, le composant stratégique, dépendant du domaine, prend les décisions conceptuelles et linguistiques; il fournit une suite ordonnée de schémas de phrases qui indique la linéarisation en phrases du texte et les items lexicaux exprimant les principaux concepts. Le deuxième, le composant syntaxique, indépendant du domaine, développe les schémas de phrases en phrases. Les modules de notre générateur sont résumés dans la Figure 1.

6. Cette phrase est inacceptable dans l'interprétation où *il* réfère à Max, interprétation qui nous intéresse ici.

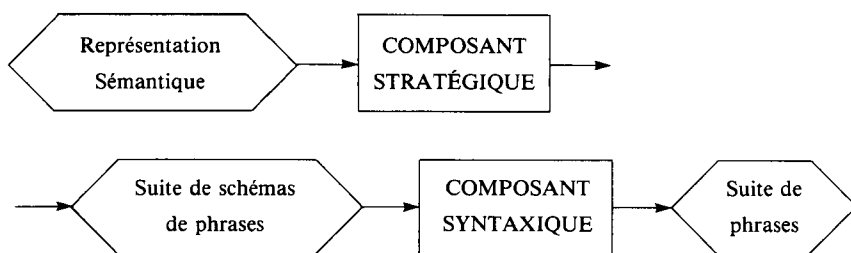


Figure 1.

Les modules de notre système de génération

La quantité d'informations syntaxiques qui doit être présente dans les schémas de phrases ainsi que le fonctionnement du composant syntaxique reposent sur les travaux du LADL. Le LADL a développé une grammaire-lexique du français couvrant 30,000 entrées utilisées dans plus de 600 types de constructions syntaxiques. De ce fait, notre composant syntaxique est basé sur une grammaire à large extension et non sur des fragments de grammaire, comme c'est trop souvent le cas dans les recherches en génération, bien que cet état de fait soit déploré par certains chercheurs (Mann, 1982).

2. Quoi dire?/Comment le dire?

Nous avons limité notre travail à la production de textes véhiculant un ensemble donné d'informations. De ce fait, nous ne traitons pas les problèmes concernant le contenu informatif

des textes produits, problèmes que nous allons brièvement évoquer.

La production d'un texte se situe dans un processus de communication: la transmission d'informations du locuteur (ici l'ordinateur) à son interlocuteur remplit une fonction, par exemple de répondre à une question sur une banque de données. La première tâche à effectuer pour répondre à une question consiste à sélectionner dans l'ensemble des données le sous-ensemble qui va constituer le contenu informatif de la réponse. Comme le souligne McKeown (1982), les informations véhiculées dans la réponse dépendent du type de question posée. Ainsi les réponses aux questions:

- Qu'est-ce que c'est que X?
- = : Qu'est-ce que c'est qu'une baleine?
- Quelle est la différence entre X et Y?
- = : Quelle est la différence entre une baleine et un marsouin?

ne contiendront pas les mêmes informations sur X: la réponse à la première question contiendra toutes les caractéristiques de X, tandis que la réponse à la seconde question ne contiendra que les caractéristiques de X qui diffèrent de celles de Y.

De plus, le contenu informatif d'un texte dépend de l'interlocuteur à qui l'on s'adresse, en particulier de ce que

l'on suppose connu de lui. Ainsi le dialogue :

Qu'est-ce que c'est qu'une sole?

Une sole est un pleuronectidé ...

se continuera par la description d'un pleuronectidé que si celle-ci est supposée inconnue de l'interlocuteur. Notons que, dans le cas d'une suite d'échanges entre le locuteur et l'interlocuteur, les informations qui ont été transmises dans les messages précédents doivent être considérées comme connues de l'interlocuteur.

La sélection des données qui constitueront le contenu informatif du texte doit se faire en dehors de toute considération linguistique. En effet, soutenir le contraire reviendrait à accepter des positions telles que: je vais exprimer la propriété P_i mais pas la propriété P_j que je ne sais pas exprimer, position inadmissible⁷.

Un modèle complet de génération de textes devrait comprendre un composant traitant de la question: quoi dire? Ce module serait en amont d'un composant traitant la question:

7. On se gardera de confondre les notions d'information non exprimée et d'information exprimée implicitement: une information non exprimée ne fait pas partie du contenu informatif du texte alors qu'une information exprimée implicitement fait partie du contenu informatif du texte, elle est inférable à partir du reste du texte.

comment le dire? Comme l'indique Mann (1982), il n'existe guère de système de génération de textes qui traite sérieusement de ces deux questions. Parmi les chercheurs qui s'intéressent au contenu informatif du texte à générer, citons d'abord Meehan (1976) qui s'occupe d'inventer des fables en s'appuyant sur les techniques de résolution de problème développées dans Schank et Abelson (1977), puis Lehnert (1981) qui se préoccupe de sélectionner et de structurer les éléments composant un résumé de récit en s'appuyant sur la trame de l'intrigue, et enfin McCoy (1981) qui sélectionne un sous-ensemble dans une banque de données pour répondre à une question, ce qui prolonge le travail de McKeown (1982).

3. Le programme qui teste le modèle

Nous avons testé notre modèle théorique de génération en réalisant un programme⁸ qui produit des textes relatant certains faits divers dans un style journalistique.

Pour illustrer la tâche accomplie par notre système, supposons qu'il existe une agence de presse internationale qui fonctionne de la manière suivante: lorsqu'un événement a lieu dans un pays, un reporter fait son enquête d'une façon traditionnelle. Ensuite, ce reporter "rédige son rapport", ce qui

8. Notre programme est écrit en LISP.

consiste à remplir un formulaire adéquat en utilisant un code indépendant de toute langue. Ce formulaire rempli est alors envoyé⁹ dans différents pays du monde où les journalistes locaux le "traduisent" dans leur langue.

Dans un tel scénario, notre système effectue le travail des journalistes qui traduisent le formulaire rempli dans la langue locale. Les formulaires remplis constituent donc nos représentations sémantiques, entrées de notre générateur. Les formulaires sont basés sur les méthodes de représentation de connaissance développées à Yale, dans le Laboratoire d'Intelligence Artificielle de R.C. Schank. Ces représentations sont fondées sur des recherches sur la mémoire (voir, entre autres, Bartlett, 1932, Minsky, 1975, Schank et Abelson, 1977 et Schank, 1982). Elles reposent sur la notion de "cadre"¹⁰, structure de mémoire représentant les connaissances spécifiques que l'on a sur un type d'événement. Les différents formulaires correspondent aux "MOPs" (Memory Organisation Packages) dans la théorie de Schank (1982). Il y a un formulaire pour les enlèvements, un pour les accidents de voiture, un pour les attentats terroristes ...

Notre système produit des textes relatant des résumés

9. Électroniquement bien sûr.

10. Nous utilisons le mot "cadre" comme traduction du mot "frame".

d'attentats terroristes commis contre des personnes. Les représentations sémantiques contiennent des informations sur l'identité de la ou des cibles (nom, nationalité, fonction...), sur l'identité du ou des terroristes (nationalité, appartenance politique...), sur les états respectifs des cibles après l'attentat (mort, blessé ou indemne), sur l'acte commis par les terroristes (coup de feu, explosion d'une bombe, embuscade...), sur la ville ou le pays où se situe l'attentat et sa date. Ces informations sont celles qui se trouvent dans un quotidien habituel et qui sont supposées intéresser le lecteur moyen. Nous utilisons ainsi implicitement un modèle de l'interlocuteur (ici le lecteur). Notons qu'une revue de la police, par exemple, supposerait un modèle différent du lecteur et ne présenterait pas le même ensemble d'informations dans les récits d'attentats.

Notre générateur fournit des textes en français et en anglais dont voici quelques exemples¹¹:

- (1) FR Un attentat a été commis contre Jacques Chirac et Bernard Pons hier à Paris: le maire de Paris a été tué mais le président du RPR est indemne. Des anarchistes ont ouvert le feu sur la voiture dans laquelle ils se rendaient à leur travail.

11. Nos exemples sont inventés. Toute ressemblance avec des personnages existant ou ayant existé est purement fortuite.

ANG Jacques Chirac and Bernard Pons were the targets of an assassination attempt yesterday in Paris: the mayor of Paris was killed but the president of RPR escaped unscathed. Anarchists opened fire on the car in which they were driving to work.

(2) FR Une bombe à retardement a explosé dans un commissariat de police aujourd'hui à Paris tuant deux policiers et en blessant dix autres. La bombe, qui était cachée dans un camion garé devant le commissariat de police, contenait vingt kilos de dynamite.

ANG A time bomb exploded in a police station today in Paris killing two policemen and wounding ten others. The bomb, which was hidden in a truck parked in front of the police contained 20 kilos of dynamite.

Le domaine choisi est intéressant dans la mesure où, relatant une suite d'événements, il met en jeu des relations temporelles, spatiales et causales. Certains travaux de génération qui portent sur des bases de données statiques (par exemple la classification des animaux ou des véhicules militaires) utilisent des solutions qui ne sont pas applicables à des données invoquant des relations temporelles, spatiales et causales. À l'inverse, les principes sur lesquels repose notre système peuvent s'appliquer à des bases de données statiques. Le fait que notre système, s'appliquant à une suite d'événements, s'applique aussi à une base de données statiques n'a rien d'étonnant dans la mesure où les relations temporelles, spatiales et locales sont associées aux phénomènes les plus ardues du langage. Signalons toutefois que la langue na-

turelle courante met pratiquement toujours en jeu de telles relations.

Il est important de souligner que le domaine des attentats, restreint d'un point de vue conceptuel, met en jeu une grande richesse linguistique tant au niveau de la syntaxe que du vocabulaire. Ce point provient du fait que la moindre variation dans la représentation sémantique peut entraîner un changement radical dans le texte produit. Cette disproportion entre variations sémantiques et syntaxiques va de pair avec l'affirmation que les rapports entre le sens et la forme sont, dans une large mesure, imprédictibles.

La capacité de notre système à produire des textes d'une grande richesse linguistique, et non des phrases plus ou moins stéréotypées et peu variées, provient du fait qu'il repose sur des principes de théorie linguistique généraux car bien vérifiés. Une conséquence de ce choix théorique est que le système est transportable d'un domaine à l'autre comme nous allons le voir.

4. Transportabilité du système

Rappelons que notre système est composé de deux modules: le composant stratégique qui prend une représentation abstraite comme entrée et fournit le schéma de phrases approprié à la représentation donnée, et le composant syntaxique qui déve-

loppe ces schémas de phrases en phrases. Seul le composant stratégique est dépendant du domaine. Nous allons donc examiner les opérations nécessaires à la construction du composant stratégique pour un domaine donné.

La première étape consiste à étudier les connaissances pragmatiques sous-jacentes au domaine et la façon de les représenter. Ceci est un point de passage obligé quel que soit le modèle de génération. Il serait en effet des plus naïfs de se lancer dans un système de production de textes sans une étude préalable du domaine. De plus, cette étude est nécessaire si les représentations abstraites servent d'autres buts que celui d'être "traduites" en langage naturel. Pour le type de représentation des connaissances que nous avons choisi, les "cadres", l'étude du domaine consiste à examiner quels sont les éléments composant un cadre donné, quelles sont les relations existant entre ces éléments et comment ils peuvent être remplis.

La deuxième étape consiste à trouver des formulations satisfaisantes quelle que soit la valeur des variables d'un cadre donné. Pour cela, il faut commencer par dégager les variables du cadre qui ont un rôle déterminant dans la synthèse. Ainsi, dans le domaine des attentats, on constate que si la cible a été spécifiquement visée, il est satisfaisant d'utiliser une formulation comme:

L'Ayatollah Khomeini a été assassiné hier à Paris.
Une bombe a explosé dans son appartement.

où le résultat de l'attentat est exprimé en premier dans une phrase indépendante. Par contre si la ou les cibles n'ont pas été visées individuellement mais en tant que représentants d'un groupe de personnes, alors une formulation où le résultat de l'attentat est exprimé au participe présent après une forme décrivant l'acte est appropriée:

- (1) Une bombe a explosé dans un commissariat de police tuant 4 policiers.

Ayant ainsi dégagé les variables pertinentes à la production et le type de formulation que l'on veut leur associer, il faut vérifier et, au besoin, ajuster le système pour obtenir des formulations acceptables pour toute la combinatoire des variables. Ainsi, il faut vérifier qu'une formulation analogue à (1) ne soit pas employée lorsque le résultat est exprimé au moyen de *assassiner*, ceci afin d'éviter la production de phrases douteuses comme:

- ?*Une bombe a explosé dans un commissariat de police assassinant 4 policiers.

Cette démarche est guidée par la connaissance des différentes façons d'exprimer une relation causale (voir par exemple Danlos, 1983), et par la connaissance des différentes constructions syntaxiques dans lesquelles entrent les groupes verbaux du lexique spécifique au domaine étudié.

Au terme de cette procédure, on a construit un ensemble de schémas de phrases avec leurs conditions d'emplois. Celles-ci sont des tests sur la représentation sémantique. Par exemple, pour un attentat, on aura des règles comme:

SI (la CIBLE est non spécifiquement visée)
 (le RÉSULTAT = MORT)
 (l'AGENT de l'attentat est non précisé)
 (l'ACTE est une EXPLOSION)

ALORS utiliser le schéma de phrases:

EXPLOSIF exploser LIEU tuant CIBLE.

La dernière étape consiste à optimiser la recherche des schémas de phrases en structurant les tests qui les conditionnent. On peut, par exemple, organiser les tests dans une structure arborescente dont les branches correspondent aux réponses "oui" ou "non" des tests.

On voit donc que la réalisation du composant stratégique pour un domaine donné demande principalement des connaissances sur le domaine lui-même et sur le lexique de ce domaine.

Jusqu'à présent, nous n'avons considéré que la production de textes, disons de la taille d'un paragraphe, "traduisant" un ensemble de données. Considérons maintenant la production de questions dans un dialogue. Tout d'abord, examinons le cas où les questions ont pour but de faire préciser

des informations manquantes ou mal comprises comme c'est le cas dans le système CADI de Nouhen-Bellec et Siroux (1981). Les types de questions sont alors conceptuellement restreints: ce sont principalement des questions en *Quel est x?* et des vérifications à la forme interrogative. Afin d'éviter de formuler toujours les vérifications dans des formes comme *Est-ce bien x?*, on peut garder un historique du dialogue. Ainsi, ayant en mémoire le nombre de vérifications déjà demandées et en particulier le nombre de fois que le même renseignement a été demandé, on est en mesure de faire varier les demandes:

Je m'excuse de vous faire répéter la même chose pour la quatrième fois, mais je n'ai toujours pas compris le département.

On notera que cet exemple n'est pas à la forme interrogative mais qu'il amènera nécessairement l'interlocuteur à répéter le nom du département. Dans ce type de dialogue, le composant stratégique sous-jacent à la formulation des questions repose-rait sur une étude de la combinatoire des types de questions et des questions déjà posées.

Dans le cas d'un dialogue où les questions ne sont pas conceptuellement aussi restreintes, par exemple:

Est-ce que vous voulez un avion qui passe par Rome, qui part à 14h et arrive à 20h ou préférez-vous un avion qui passe par Nice, part à 10h et arrive à 17h?

le composant stratégique produisant les schémas de phrases de

telles questions complexes doit avoir connaissance de la combinatoire des renseignements qui peuvent être demandés et des différentes façons de les articuler. Il serait alors analogue aux composants stratégiques construits pour "traduire" un ensemble de données. Bien entendu, il peut être complété par un historique des questions déjà posées afin de formuler les demandes de vérification comme dans le cas précédent.

En conclusion, notre modèle de génération, conçu dans l'optique de générer des paragraphes, peut aisément être appliqué pour "traduire" - en plusieurs langues - un ensemble de données, que celui-ci concerne les attentats terroristes (domaine qui a servi de test à notre modèle) ou les conditions météorologiques, par exemple. De plus, il semble que ce modèle soit facilement adaptable pour produire des questions (conceptuellement restreintes ou non) dans le cadre d'un dialogue.

5. Langage écrit/parlé

Notre système, générant des textes dans le style journalistique, simulerait un locuteur écrivant plutôt qu'un locuteur parlant. En particulier, notre générateur produit des textes bien construits, ce qui n'est pas toujours le type de textes formulés dans un discours oral.

Bon nombre de travaux en génération automatique tendent

à simuler un locuteur parlant (voir, entre autres, McGuire (1980), Hovy (1983), McDonald (1982)). Ces générateurs, qui se veulent "plausibles psychologiquement", s'appuient sur des observations comme: lorsqu'une personne parle, elle ne sait pas exactement quelles idées elle va exprimer. Ou bien, lorsque l'on commence une phrase, on ne sait pas à l'avance les mots que l'on va employer. De ce fait, l'idée de base de ces générateurs est de produire le texte de "gauche à droite". Ainsi, pour des phrases affirmatives de langues à ordre de base sujet-verbe-complément, la procédure consiste à déterminer un terme qui va constituer le sujet et le synthétiser; ensuite, le mécanisme recherche un terme pour le verbe et le synthétise, et ainsi de suite. Si les observations sur lesquelles se basent ces générateurs ne sont pas sans quelque fondement, elles demanderaient à être sensiblement précisées et raffinées, car l'hypothèse que l'on produise un texte pas à pas, de "gauche à droite", sans anticiper sur ce que va être la suite, nous paraît grossièrement fausse. Certes, si les discours parlés sont moins soignés que les textes écrits, ce ne sont quand même jamais des magmas informes et incohérents tels que devraient l'être des textes produits de "gauche à droite". De plus, lorsqu'un texte est effectivement incohérent du fait d'un lapsus:

J'ai laissé ma valise dans mon cigare.

ces lapsus montrent généralement que le locuteur anticipe sur ce qu'il allait dire. L'étude des lapsus faite par Garrett (1980) ainsi que d'autres travaux de psychologues (voir, entre autres, Miller *et alii* (1960), Goldman-Eisler (1980) et Butterworth (1980), vont dans le sens de nos propositions, à savoir que le locuteur planifie à l'avance ce qu'il va dire. Ces mêmes travaux de psychologues montrent aussi que le locuteur se souvient de ce qu'il a déjà dit et que ce qu'il est en train d'exprimer à un moment donné dépend de la séquence qui précède comme de celle qui suit. Dans les générateurs de "gauche à droite", toutes les décisions sont locales: elles ne concernent que ce qui est en train d'être exprimé. Non seulement l'application de décisions locales est susceptible de fournir des textes incohérents, mais plus grave encore, ces textes risquent d'être agrammaticaux. Ainsi, le composant syntaxique du générateur de McDonald (1982), qui ne repose que sur des décisions locales, peut tenir compte des dépendances de gauche à droite (par l'intermédiaire de registres) mais pas des dépendances de droite à gauche, ni des phénomènes qui demandent une vue d'ensemble de la phrase. McDonald justifie cet état de fait en s'appuyant sur l'hypothèse que seules les dépendances syntaxiques séquentielles peuvent être appréciées lors d'un discours oral. De nouveau, cette hypothèse nous paraît grossièrement fautive. À titre d'exemple, considérons la paire suivante:

(1) a. Quand Max chante, il agace Ève.

(2) a. Quand il chante, Max agace Ève.

Il semble qu'il ne soit pas plus difficile de produire (2) que (1). Or, dans (1), on a une dépendance de gauche à droite (le pronom est placé à droite du nom auquel il réfère), tandis que dans (2), on a une dépendance de droite à gauche (le pronom est placé à gauche du nom auquel il réfère). La situation est analogue avec des paires mettant en jeu la réduction de complétive à l'infinitive:

(1) b. Max aime chanter.

(2) b. Chanter plaît à Max.

On peut donc s'interroger sur la valeur d'un générateur qui peut produire (1) naturellement mais qui ne peut produire (2) que d'une façon *ad hoc* ou "*post hoc*" (sic) (McDonald, 1982, p. 26).

Un des principes de base de notre modèle de génération est que la production d'un texte repose sur des décisions non locales: le problème de la linéarisation en phrases et celui du choix adéquat des mots demandent tous deux que la représentation sémantique soit considérée dans son ensemble. Un texte doit s'appuyer sur une structure de discours globale qui impose ses contraintes lors de la synthèse des différents éléments composants. Ces structures de discours sont mises en correspondance avec des classes de représentations sémantiques sur

la base de propriétés éparpillées dans les représentations sémantiques.

Nous ne prétendons pas avoir construit un modèle de génération qui simule un locuteur écrivant. Ne serait-ce que parce que notre système fournit directement le texte définitif sans passer par des étapes de corrections, points de passage obligés de tout rédacteur. Mais notre modèle, basé sur des structures de discours où les décisions conceptuelles et linguistiques sont prises simultanément, est certainement plus proche de la réalité que les modèles de "gauche à droite". Enfin et surtout, il permet de produire (en plusieurs langues) des textes bien construits¹², d'une grande variété linguistique et s'appliquant à des domaines invoquant des relations temporelles, spatiales et causales.

Laurence Danlos

*Laboratoire d'automatique
documentaire et linguistique
(C.N.R.S.), Paris*

12. Signalons, à ce propos, que si l'on s'intéresse à des applications pratiques d'un système de génération automatique, il est nécessaire de fournir des textes de bonne qualité: quel serait l'utilisateur qui voudrait d'un système de questions/réponses dont les réponses seraient incomplètes, mal construites et d'un style relâché?

RÉFÉRENCES

- BARTLETT, F.C. (1932) *Remembering, a study in experimental and social psychology*, Cambridge (England), Cambridge University Press.
- BUTTERWORTH, B. (1980) "Evidence from Pauses in Speech", dans *Language Production*, London, B. Butterworth ed., Academic Press.
- DANLOS, L. (1981) "La morphosyntaxe des expressions figées", *Langages*, n° 63, Paris, Larousse.
- DANLOS, L. (1983) "Génération automatique de textes en langues naturelles", thèse d'État, Université de Paris 7.
- GARRETT, M.F. (1980) "Levels of Processing in Sentence Production" dans *Language Production*, London, B. Butterworth ed., Academic Press.
- GIRY-SCHNEIDER, J. (1981) "Les compléments nominaux du verbe DIRE", *Langages*, n° 63, Paris, Larousse.
- GREEN, G.M. (1972) "Some observations on the syntax and semantics of instrumental verbs" dans *Papers from the eight regional meeting Chicago Linguistic Society*, Chicago, Chicago Linguistic Society.
- GUILLET, G. et C. LECLÈRE (1981) "Restructuration du groupe nominal", *Langages*, n° 63, Paris, Larousse.
- GOLDMAN-EISLER, F. (1980) "Psychological Mechanisms of Speech Production as studied through the Analysis of Simultaneous Production" dans *Language Production*, London, B. Butterworth ed., Academic Press.
- GROSS, M. (1981) "Les bases empiriques de la notion de prédicat sémantique", *Langages*, n° 63, Paris, Larousse.
- HARRIS, Z. (1982) *A Grammar of English on Mathematical Principles*, New York, Wiley-Interscience.
- HOVY, E. (1983) "Computer-generated Language", Yale University (inédit).
- LEHNERT, W.G. (1981) "Plot Units and Narrative Summarization", *Cognitive Science*, vol. 4.
- MANN, W. (1982) "Text Generation", *American Journal of Computational Linguistics*, vol. 8, n° 2.
- MCCAWLEY, J.D. (1971) "Prelexical Syntax" dans *Report of the 22nd annual round table meeting on Linguistics and Language Studies*, O'Brien ed., Georgetown University Press.
- MCCOY, K.F. (1981) "Automatic Enhancement of a Database Knowledge Representation for Natural Language Generation", Technical Report MS-CIS-81-6, University of Pennsylvania.
- MCDONALD, D. (1982) "Natural Language Generation as a Computational Problem: an introduction", dans *Computational Theories of Discourse*, Cambridge (Mass.), Brady ed., MIT Press.

- McGUIRE, R. (1980) "Political primaries and words of pain", Yale University (inédit).
- McKEOWN, K.R. (1982) *Generating Natural Language Text in response to Questions about database structure*, thèse de Ph. D., University of Pennsylvania.
- MEEHAN, J. (1976) *The metanovel: Writing stories by computer*, thèse de Ph. D., Yale University.
- MILLER, G.A., E. GALANTER et K.H. PRIBAM (1960) *Plans and the Structure of Behavior*, New York, Holt, Rinehart and Winston Inc.
- MINSKY, M. (1975) "A framework for representing knowledge" dans *The psychology of computer vision*, New York, P.H. Winston ed., McGraw-Hill.
- NOUHEN-BELLECC, A. et J. SIROUX (1981) *CADI: constructeur automatique de dialogue intelligent intégré à un système de reconnaissance et de compréhension de la parole*, thèse de 3e cycle, Université de Rennes.
- SCHANK, R.C. (1982) *Dynamic memory: A theory of learning in computers and people*, Cambridge (Mass.), Cambridge University Press.
- SCHANK, R.C. et R.P. ABELSON (1977) *Scripts, plans, goals, and understanding*, Hillsdale (N.Jersey), Lawrence Erlbaum Associates.