

# Génération de réponses en langue naturelle orales et écrites pour les systèmes de question-réponse en domaine ouvert

**Anne Garcia-Fernandez**  
Soutenance de thèse

LIMSI-CNRS et Université Paris Sud 11 Orsay

10 décembre 2010

## Contexte 1/2

Qu'est-ce que répondre à une question en domaine ouvert ?

- Des **systèmes** ont pour but de **donner une réponse** précise à **une question**.
- Ils **extraient** une ou plusieurs réponses à **partir de documents**.
- Ils sont **évalués** lors de **campagnes** [Voorhees, 1999] [Galibert, 2009]

[Q] Où est-ce que se trouve la Joconde ?

[R] au Louvre

### Notre point de vue

- Réponse à une question vue comme une **interaction**
- Fournir une **réponse en langue naturelle**
- Ne se limitant pas à donner l'information qui répond à la question

[R] La Joconde se trouve au Louvre.

## Contexte 2/2

### Notre problématique

Quelle **forme linguistique** choisir ?

[R] au Louvre

[R] La Joconde se trouve au Louvre.

[R] Elle est au Louvre.

[R] C'est au Louvre que la Joconde se trouve.

### Travaux traitant ce point

- des travaux en domaine fermé [*Benamara, 2003*]
- objectifs différents [*Plamondon, 2002*]

### Première idée

Chercher **un corpus**

## Utiliser un corpus de réponses existant ?

Réponses des systèmes lors des campagnes d'évaluation

[R] au Louvre, à Paris, en France, ...

☹ **Réponses succinctes**                      ☹ **Non spontanées**

☹ **Trop limitées**

Réponses issues de sites collaboratifs

[R] La Joconde de Leonard de Vinci, qui se trouve au Louvre depuis 1804, peut être admirée dans une salle spécialement aménagée pour la recevoir, la salle des États. (...)

☹ **Réponses complexes**                      ☹ **Présentant trop d'informations**

☹ **Corpus hétérogène**

### Deuxième idée

Chercher un **corpus de questions** et les **poser à des locuteurs**.

# Constituer un corpus de réponses d'humains

Collecter un corpus de réponses

Qui puissent être produites par des systèmes

Analyser ces réponses

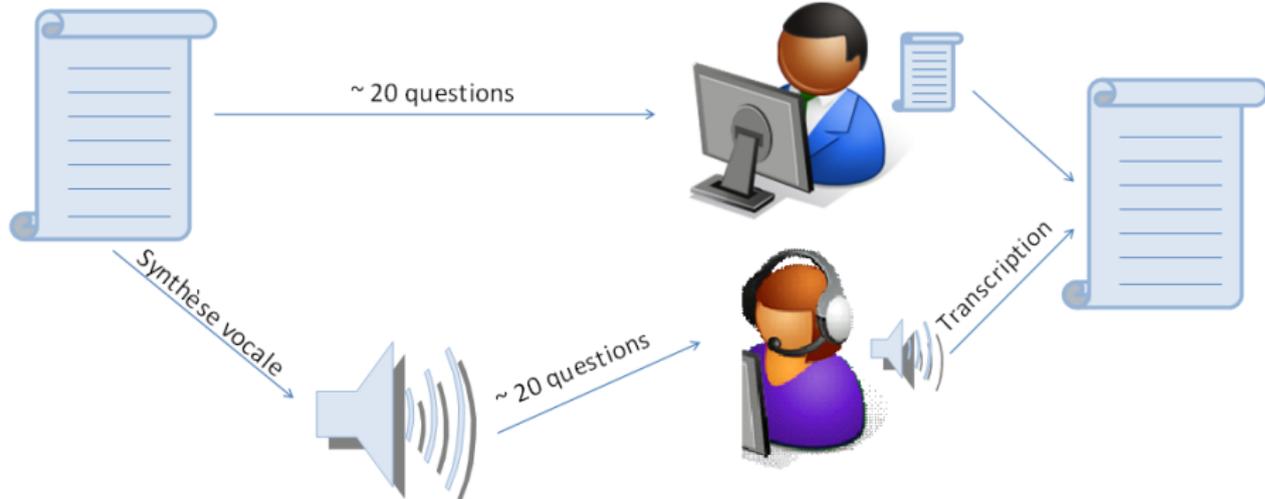
Pour en extraire des *règles* sur lesquelles un module de génération de réponse en langue naturelle pourrait s'appuyer

# Plan

- 1 Introduction
- 2 Protocole expérimental
- 3 Construction du corpus de questions
- 4 Collecte d'un corpus de réponses
- 5 Analyses du corpus
- 6 Conclusion et perspectives

# Protocole expérimental

Des questions construites



## Scénario (ou contexte expérimental)

“ La classe de CM2 de l'école primaire de Château-la-Vallière se propose de créer des affiches pour la fête de l'école. Ces posters portent sur des sujets variés et pour les réaliser, les élèves ont besoin de quelques informations. Vous êtes là pour répondre à leurs questions !  
Ne soyez pas surpris s'il vous semble avoir déjà répondu à une question. Certaines se ressemblent en effet mais ont été posées par différents enfants. ”

# Objectifs du protocole

## Obtenir des réponses. . .

- spontanées

## Contraintes sur le protocole :

- répondre à des enfants de CM2

# Objectifs du protocole

## Obtenir des réponses. . .

- spontanées
- à valence positive

## Contraintes sur le protocole :

- répondre à des enfants de CM2
- questions “*faciles*”

# Objectifs du protocole

## Obtenir des réponses. . .

- spontanées
- à valence positive
- non succinctes

## Contraintes sur le protocole :

- répondre à des enfants de CM2
- questions “*faciles*”
- répondre à des enfants de CM2

# Objectifs du protocole

## Obtenir des réponses. . .

- spontanées
- à valence positive
- non succinctes
- reproductibles par des systèmes

## Contraintes sur le protocole :

- répondre à des enfants de CM2
- questions “*faciles*”
- répondre à des enfants de CM2
- questions précises

# Objectifs du protocole

## Obtenir des réponses. . .

- spontanées
- à valence positive
- non succinctes
- reproductibles par des systèmes
- comparables à l'oral et à l'écrit

## Contraintes sur le protocole :

- répondre à des enfants de CM2
- questions “*faciles*”
- répondre à des enfants de CM2
- questions précises
- deux protocoles très similaires

# Objectifs du protocole

## Obtenir des réponses. . .

- spontanées
- à valence positive
- non succinctes
- reproductibles par des systèmes
- comparables à l'oral et à l'écrit
- non "systématiques"

## Contraintes sur le protocole :

- répondre à des enfants de CM2
- questions "*faciles*"
- répondre à des enfants de CM2
- questions précises
- deux protocoles très similaires
- nombre et ordre des questions pour un participant

# Plan

- 1 Introduction
- 2 Protocole expérimental
- 3 Construction du corpus de questions**
- 4 Collecte d'un corpus de réponses
- 5 Analyses du corpus
- 6 Conclusion et perspectives

## De l'importance des questions posées

La GRammaire INTeractive, une théorie de la réponse en interaction  
 Met en valeur le lien entre **la forme linguistique d'une question** et **la réponse présupposée** par le locuteur questionnant [*Luzzati, 2001*]

[Q] Je voudrais savoir où est la Joconde?    [R] euh... je sais pas  
 [Q] La Joconde est où?                            [R] au Louvre

### Contraintes

- Corpus contenant **différentes variations** autour d'une même question
- Des questions de **forme linguistique contrôlée**

# Questions de base

## Principe

Variations linguistiques à partir de **questions de base** choisies pour :

- leur thème
- leur type sémantique
- leur réponse attendue
- le type sémantique de leur réponse attendue

# Questions de base

## Principe

Variations linguistiques à partir de **questions de base** choisies pour :

- **leur thème : culture générale, questions “faciles”**
- leur type sémantique
- leur réponse attendue
- le type sémantique de leur réponse attendue

(theme<sub>1</sub>) Lieu d'exposition de la Joconde

(theme<sub>2</sub>) Durée d'un mandat présidentiel

(theme<sub>3</sub>) Taille d'un double décimètre

...

# Questions de base

## Principe

Variations linguistiques à partir de **questions de base** choisies pour :

- leur thème
- **leur type sémantique : quantité, lieu, temps [Lailier, 2009]**
- leur réponse attendue
- le type sémantique de leur réponse attendue

(lieu) Où est la Joconde ?

(temps) Quand ont lieu les Jeux Olympiques d'été ?

(quantité) Combien mesure un double décimètre ?

# Questions de base

## Principe

Variations linguistiques à partir de **questions de base** choisies pour :

- leur thème
- leur type sémantique
- **leur réponse attendue : unique ou multiple**
- le type sémantique de leur réponse attendue

(unique) Combien mesure un double décimètre ?

(multiple) Combien pèse un bébé à la naissance ?

# Questions de base

## Principe

Variations linguistiques à partir de **questions de base** choisies pour :

- leur thème
- leur type sémantique
- leur réponse attendue
- **le type sémantique de leur réponse attendue : spécifique au type sémantique de la question**

(taille) Combien mesure un double décimètre ?

(poids) Combien pèse une brique de lait ?

(durée) Combien dure un mois de février ?

# Questions de base

Au total, 19 questions de base

Quantité	<p>Combien dure un mois de février ?</p> <p>Combien dure une grossesse ?</p> <p>Combien dure un mandat présidentiel en France ?</p> <p>Combien mesure un bébé à la naissance ?</p> <p>...</p>
Lieu	<p>Où arrivent les avions à Paris ?</p> <p>Où arrive le Tour de France ?</p> <p>Où se trouve la Joconde ?</p> <p>Où se trouvent les Alpes ?</p> <p>...</p>
Temps	<p>Quand a eu lieu l'armistice de la 1ère guerre mondiale ?</p> <p>Quand ont lieu les jeux olympiques d'été ?</p>

## Variations des questions 1/2

### Variation morphosyntaxique

La GRINT modélise les variations des questions selon deux dimensions :

- syntagmatique : variation de la forme syntaxique
- paradigmatique : variation de l'interrogatif

### Exemples de variations syntagmatiques

Où est la Joconde ?

La Joconde est où ?

Je voudrais savoir où est la Joconde ?

### Exemples de variations paradigmatiques

Où est la Joconde ?

Dans quel musée est la Joconde ?

La Joconde est-elle au Louvre ?

## Variations des questions 2/2

### Variation du verbe principal

- verbe spécifique
- verbe neutre

(spécifique) Quand **ont lieu** les jeux olympiques d'été ?

(neutre) Quand **sont** les jeux olympiques d'été ?

### Granularité

- neutre : pas de type précis demandé
- fine : type plus précis que ↓
- grosse : type plus large que ↑

(neutre) Dans quel **endroit** se trouve la Joconde ?

(fin) Dans quel **musée** se trouve la Joconde ?

(gros) Dans quelle **ville** se trouve la Joconde ?

## Le corpus de questions en quelques chiffres

- 707 questions : 207 de quantité, 400 de lieu et 100 de temps
- plus de 35 variations par question de base
- des questions ouvertes et fermées
- formatées pour l'écrit (majuscules, . . .) et synthétisées pour l'oral

### Nombre de participants nécessaires à la collecte

#### Contraintes

- au moins 3 réponses par variation de question
- ne pas poser *trop* de questions à chaque participant

Au minimum 144 participants ( $\frac{1}{2}$  à l'oral,  $\frac{1}{2}$  à l'écrit) pour des séries de 18 à 24 questions

# Plan

- 1 Introduction
- 2 Protocole expérimental
- 3 Construction du corpus de questions
- 4 Collecte d'un corpus de réponses**
- 5 Analyses du corpus
- 6 Conclusion et perspectives

# Déroulement des passations

## Recrutement des participants

Personnes contactées	2164
Personnes ayant répondu	12 %
Personnes ayant accédé à l'une des plates-formes	95 %
Passations effectuées	63 %
Nombre de passations valides	152

- Soit **7 % des personnes contactées** ont effectué une **participation valide** à l'expérience

## Le corpus de réponses en quelques chiffres

	Tout	Oral	Écrit
Nombre de réponses	<b>3132</b>	1044	2088
Nombre de questions différentes	<b>707</b>	493	707
Nombre de participants	152	<b>53</b>	<b>99</b>
Nombre de réponses/question	<b>6,17</b>	2,39	4,11

- Durée du corpus oral : ~ 1h
- Nombre total de mots : 25 104
- Nombre moyen de mots par réponse : 8,01

## Exemples de réponses collectées

### Réponses à des questions ouvertes

**A31** : Au Louvre

**A1938** : un mois de février dure à peu près 28 jours

**A60** : 28 jours les années bissextiles, 29 jours les autres

### Réponses à des questions fermées

**A408** : Tout à fait !

**A1404** : oui, 28 jours les années normales, 29 les bissextiles

**A2499** : Elle a bien été signée à cette date

### Réponses à valence négative

**A2480** : je n'en ai aucune idée

### Réponses complexes et complétives

**A2600'** : Les avions peuvent atterrir sur différents aéroports autour de Paris. Globalement ils se répartissent ainsi : Charles de Gaulle à Roissy pour (...)

# Évaluation de la collecte 1/4

## Spontanéité des réponses

### 😊 Temps de réponse

- à l'oral : moins de 7 secondes
- à l'écrit : entre 40 sec et 2 min

### 😊 Indices qualitatifs

- à l'oral : présence d'hésitations
- en général : présence d'expressions de doute
- en général : humour, abréviations, ...

**A1138** : euh oui c'est à dire que euh la section graduée mesure 20 centimètres

**A521** : Dans plusieurs, l'Allemagne, la France et ptetre la Belgique

**A1582** : Si seulement on pouvait se débarrasser de cette crapule plus tôt...

# Évaluation de la collecte 2/4

## Définition : information-réponse

Nouvelle information qui correspond soit au type attendu de la réponse, soit à un aveu d'incompétence ("je ne sais pas", par exemple)

**A442'** : Une brique de lait pèse autour de 6 kilos

**A1262** : je n'en sais rien [rire]

Facilité de fournir une réponse aux questions

😊 Presque 95 % des réponses contiennent une **information-réponse**

Des réponses non succinctes

😞😊 57 % des réponses sont constituées d'une **information-réponse** seule

# Évaluation de la collecte 3/4

## Définition : information complémentaire

Élément qui apporte une nouvelle information qui ne correspond pas à l'information-réponse

**A2904** : 20,9 cm en largeur, 29,7 cm en hauteur, 0,23 mm en profondeur. D'une blancheur immaculée.

## Des réponses sans information complémentaire

😊 Moins de 10 % des réponses contiennent au moins une information complémentaire

# Évaluation de la collecte 4/4

Des réponses que l'on peut comparer

😊 Entre oral et écrit :

- 1044 réponses à l'oral, 2088 à l'écrit
- obtenues par des protocoles très similaires

😊 Pour une même question de base :

- au moins 130 réponses par "groupe de variation"

😊 Pour chaque variation d'une question de base

- au moins 6 réponses par variation de question

# Plan

- 1 Introduction
- 2 Protocole expérimental
- 3 Construction du corpus de questions
- 4 Collecte d'un corpus de réponses
- 5 Analyses du corpus**
- 6 Conclusion et perspectives

# Analyse du corpus

## Questions à se poser

- ❶ Est-il judicieux de réutiliser des éléments de la question ?
- ❷ Quelle information faut-il répondre ?
- ❸ Comment construire une réponse minimale ?
- ❹ La modalité (oral/écrit) influe-t-elle sur la formulation de réponse à générer ?
- ❺ Euh...et si on hésite ?
- ❻ Que répondre quand on n'a pas de réponse ?
- ❼ Comment répondre plusieurs informations-réponse ?

# 1. Est-il judicieux de réutiliser des éléments de la question ?

## Observation des reprises de la question

- Annotation des différents éléments des questions
- Annotation de ces éléments dans les réponses

verbe objet type info. complémentaire réponse-suggérée

Q258 : Où est la Joconde ?

A2663 : La Joconde est actuellement au Louvre

Q006 : Combien de kilos pèse un bébé à la naissance ?

A676 : un bébé à la naissance pèse environ 3 kilos 5

Q212 : La Joconde se trouve au Louvre ?

A2389 : La Joconde est en effet actuellement exposée au Louvre.

# 1. Est-il judicieux de réutiliser des éléments de la question ?

## Annotation de cas particuliers

- Abréviations (des types notamment)
- Formes erronées (fautes typographiques, orthographiques)
- Formes fléchies (singulier/pluriel, modalité. . .)

verbe   objet   type   info. complémentaire   réponse-suggérée

A45' : Une briquet de 1L doit peser environ 1 kg (...)

Une briquet au lieu de "brique"

doit peser au lieu de "pèse"

kg au lieu de "kilo"

# 1. Est-il judicieux de réutiliser des éléments de la question ?

Reprise de l'objet : 3 types

- Reprise exacte
- Reprise avec modification
- Reprise par un pronom

verbe    objet    type    info. complémentaire    réponse-suggérée

**Q049** : Combien est-ce que **pèse** **une bouteille d'eau** ?

**A2280** : **Une bouteille d'eau** contient du liquide. (...) Si **la bouteille** contient 1 litre, **elle** **pèsera** un kilo et ainsi de suite.

**Une bouteille d'eau** : reprise exacte

**la bouteille** : reprise avec modification

**elle** : reprise par un pronom

# 1. Est-il judicieux de réutiliser des éléments de la question ?

Élément repris (%)	Tout	Oral	Écrit
Tous	<b>9.90</b>	11.69	9.02
$\geq 1$ élément	<b>27.39</b>	29.23	26.50
Objet	<b>22.54</b>	22.31	22.65
Verbe	15.79	17.83	14.78
Type	14.51	<b>19.20</b>	<b>12.08</b>
Information complémentaire	11.79	12.59	11.27
Réponse suggérée	9.58	13.26	7.87

# 1. Est-il judicieux de réutiliser des éléments de la question ?

Élément repris (%)	Tout	Oral	Écrit
Tous	<b>9.90</b>	11.69	9.02
>= 1 élément	<b>27.39</b>	29.23	26.50
Objet	<b>22.54</b>	22.31	22.65
Verbe	15.79	17.83	14.78
Type	14.51	<b>19.20</b>	<b>12.08</b>
Information complémentaire	11.79	12.59	11.27
Réponse suggérée	9.58	13.26	7.87

Parmi les réponses qui reprennent l'objet

	Tout	Oral	Écrit
Ont au moins une reprise exacte	<b>62.48</b>	67.24	60.16
Ont au moins une reprise avec modification	<b>16.11</b>	14.41	16.94
Ne le reprennent que par un pronom	23.39	<b>18.77</b>	<b>25.63</b>

## 2. Quelle information faut-il donner en réponse ?

### Annotation de l'information-réponse

Comme l'élément de la réponse qui :

- soit est une nouvelle information qui correspond au type attendu par la question (valence positive)
- soit exprime un aveu d'incompétence (valence négative)

### Réponse à "valence positive"

**A1150** : euh un mandat présidentiel dure 5 ans en France

**A783** : Oui, puisque 1 décimètre c'est 10 centimètres.

### Réponse à "valence négative"

**A793** : J'avoue que je ne sais pas du tout. (...)

## 2. Quelle information faut-il donner en réponse ?

### Observation des annotations

3 grands types d'informations-réponse :

- à valence positive, à une question ouverte
- à valence positive, à une question fermée
- à valence négative

Des réponses :

- sans information-réponse (3 % du corpus)
- avec plusieurs informations-réponse (19 %)

## 2. Quelle information faut-il donner en réponse ?

### Observation des annotations

3 grands types d'informations-réponse :

- à **valence positive**, à **une question ouverte**
- à valence positive, à une question fermée
- à valence négative

Des réponses :

- sans information-réponse (3 % du corpus)
- avec plusieurs informations-réponse (19 %)

### Réponse à valence positive à une question ouverte

**A1150** : euh un mandat présidentiel dure 5 ans en France

**A2488** : Le Tour de France fini par Paris

**A2937** : Elle a eu lieu le 11 novembre 1918.

## 2. Quelle information faut-il donner en réponse ?

### Observation des annotations

3 grands types d'informations-réponse :

- à valence positive, à une question ouverte
- à **valence positive, à une question fermée**
- à valence négative

Des réponses :

- sans information-réponse (3 % du corpus)
- avec plusieurs informations-réponse (19 %)

### Réponses à valence positive à des questions fermées

**A408** : Tout à fait !

**A1404** : oui, 28 jours les années normales, 29 les bissextiles

**A2499** : Elle a bien été signée à cette date

**A2608** : Le Rhône ne finit pas dans le delta de la Camargue

## 2. Quelle information faut-il donner en réponse ?

### Observation des annotations

3 grands types d'informations-réponse :

- à valence positive, à une question ouverte
- à valence positive, à une question fermée
- à **valence négative**

Des réponses :

- sans information-réponse (3 % du corpus)
- avec plusieurs informations-réponse (19 %)

### Réponse à valence négative

**A793** : J'avoue que je ne sais pas du tout. (...)

**A2480** : je n'en ai aucune idée

## 2. Quelle information faut-il donner en réponse ?

### Observation des annotations

3 grands types d'informations-réponse :

- à valence positive, à une question ouverte
- à valence positive, à une question fermée
- à valence négative

### Des réponses :

- **sans information-réponse (3 % du corpus)**
- **avec plusieurs informations-réponse (19 %)**

### 0 à plusieurs réponses

**A1121** : j' ai pas compris la question

**A60** : 28 jours les années bissextiles, 29 jours les autres

**A2525'** : au Louvre, à Paris, en France

## 2. Quelle information faut-il donner en réponse ?

### Observation des annotations

3 grands types d'informations-réponse :

- à valence positive, à une question ouverte
- à valence positive, à une question fermée
- à valence négative

Des réponses :

- sans information-réponse (3 % du corpus)
- avec plusieurs informations-réponse (19 %)

Comment l'information-réponse et les reprises de la question s'organisent-elles au sein d'une réponse ?

### 3. Comment construire une réponse minimale ?

#### Idée

Ne tenir compte que des éléments dont un système de question-réponse peut disposer

#### Méthode

Construction de patrons de réponse

Forme de la réponse réduite à :

- l'information-réponse
- les éléments de la question repris

**A45** : Une brique de 1L doit peser environ 1 kg (...)

⇒ **Patron** : objet verbe info-rep type

### 3. Comment construire une réponse minimale ?

#### Patrons les plus fréquents

	Tout	Oral	Écrit
- info-rep -	2100	711	1389
objet verbe info-rep -	147	54	93
objet_pronom verbe info-rep -	56	20	36
objet info-rep -	44	11	33
objet_modifié verbe info-rep -	23	8	15
objet verbe info-rep info-complémentaire	22	11	11

#### Synthèse

- Ordre des éléments figé : objet verbe info-rep
- Information-réponse après les éléments repris de la question
- Patrons trop limités
- Du temps exploitable par les systèmes

## 4. La modalité influence-t-elle sur la formulation de la réponse ?

### Points communs

- Éléments de la question les plus repris
- Forme des informations-réponse
- Patrons de réponse les plus fréquents
- Nombre d'informations-réponse proposées
- Nombre de réponses succinctes

## 4. La modalité influe-t-elle sur la formulation de la réponse ?

## Différences

	Oral	Écrit
Réponses complétives	5 %	10 %
Reprises par pronom	19 %	26 %
Temps de réponse	<10 sec	1 à 2 min
Longueur moyenne des réponses	6 mots	8 mots
Lexique en nombre de formes distinctes	1 634	3 363
Taux de réponses contenant plusieurs informations-réponse	14 %	22 %

## 4. La modalité influe-t-elle sur la formulation de la réponse ?

## Différences

	Oral	Écrit
<b>Réponses complétives</b>	5 %	10 %
Reprises par pronom	19 %	26 %
Temps de réponse	<10 sec	1 à 2 min
Longueur moyenne des réponses	6 mots	8 mots
Lexique en nombre de formes distinctes	1 634	3 363
Taux de réponses contenant plusieurs informations-réponse	14 %	22 %

**A878** : 5 ans depuis 2002, avant il durait 7 ans. C'est l'ancien président Jacques Chirac qui a fait modifier cette durée.

## 4. La modalité influe-t-elle sur la formulation de la réponse ?

## Différences

	Oral	Écrit
Réponses complétives	5 %	10 %
<b>Reprises par pronom</b>	19 %	26 %
Temps de réponse	<10 sec	1 à 2 min
Longueur moyenne des réponses	6 mots	8 mots
Lexique en nombre de formes distinctes	1 634	3 363
Taux de réponses contenant plusieurs informations-réponse	14 %	22 %

**A657** : il dure 28 jours euh les années non bissextiles et 29 jours les années bissextiles

## 4. La modalité influe-t-elle sur la formulation de la réponse ?

## Différences

	Oral	Écrit
Réponses complétives	5 %	10 %
Reprises par pronom	19 %	26 %
<b>Temps de réponse</b> ☺	<10 sec	1 à 2 min
Longueur moyenne des réponses	6 mots	8 mots
Lexique en nombre de formes distinctes	1 634	3 363
Taux de réponses contenant plusieurs informations-réponse	14 %	22 %

Mesure différente selon la modalité

## 4. La modalité influe-t-elle sur la formulation de la réponse ?

## Différences

	Oral	Écrit
Réponses complétives	5 %	10 %
Reprises par pronom	19 %	26 %
Temps de réponse ☹	<10 sec	1 à 2 min
<b>Longueur moyenne des réponses</b>	6 mots	8 mots
Lexique en nombre de formes distinctes	1 634	3 363
Taux de réponses contenant plusieurs informations-réponse	14 %	22 %

**A758** : Une brique de lait contient 1L de lait. Or, 1L pèse 1Kg. Donc une brique de lait pèse 1Kg.

**A1225** : une brique de lait pèse euh pèse un kilo

## 4. La modalité influe-t-elle sur la formulation de la réponse ?

## Différences

	Oral	Écrit
Réponses complétives	5 %	10 %
Reprises par pronom	19 %	26 %
Temps de réponse ☹	<10 sec	1 à 2 min
Longueur moyenne des réponses	6 mots	8 mots
<b>Lexique en nombre de formes distinctes ☹</b>	1 634	3 363
Taux de réponses contenant plusieurs informations-réponse	14 %	22 %

- Pas de correction sur le sous-corpus de réponses écrites
- Deux fois plus de réponses à l'écrit qu'à l'oral

## 4. La modalité influe-t-elle sur la formulation de la réponse ?

## Différences

	Oral	Écrit
Réponses complétives	5 %	10 %
Reprises par pronom	19 %	26 %
Temps de réponse ☹	<10 sec	1 à 2 min
Longueur moyenne des réponses	6 mots	8 mots
Lexique en nombre de formes distinctes ☹	1 634	3 363
<b>Taux de réponses contenant plusieurs informations-réponse</b>	14 %	22 %

**A2525'** : au Louvre, à Paris, en France

**A2479** : Oui, effectivement

## Analyse du corpus : Synthèse

- 1 **Est-il judicieux de réutiliser des éléments de la question ?**
- 2 Quelle information faut-il répondre ?
- 3 Comment construire une réponse minimale ?
- 4 La modalité (oral/écrit) influe-t-elle sur la formulation de réponse à générer ?
- 5 Euh...et si on hésite ?
- 6 Que répondre quand on n'a pas de réponse ?
- 7 Comment répondre plusieurs informations-réponse ?

- 27 % des réponses reprennent au moins l'un des éléments de la question
- Objet, verbe puis type sont les éléments les plus repris
- Trois types de reprise de l'objet : la reprise exacte est la plus courante

## Analyse du corpus : Synthèse

- 1 Est-il judicieux de réutiliser des éléments de la question ?
- 2 **Quelle information faut-il répondre ?**
- 3 Comment construire une réponse minimale ?
- 4 La modalité (oral/écrit) influe-t-elle sur la formulation de réponse à générer ?
- 5 Euh...et si on hésite ?
- 6 Que répondre quand on n'a pas de réponse ?
- 7 Comment répondre plusieurs informations-réponse ?

3 grands types d'informations-réponse :

- en fonction du type ouvert ou fermé de la question
- en fonction de la valence de la réponse

19% de réponses contenant plusieurs informations-réponse

## Analyse du corpus : Synthèse

- ❶ Est-il judicieux de réutiliser des éléments de la question ?
  - ❷ Quelle information faut-il répondre ?
  - ❸ **Comment construire une réponse minimale ?**
  - ❹ La modalité (oral/écrit) influe-t-elle sur la formulation de réponse à générer ?
  - ❺ Euh...et si on hésite ?
  - ❻ Que répondre quand on n'a pas de réponse ?
  - ❼ Comment répondre plusieurs informations-réponse ?
- 
- Ordre des éléments figé : objet verbe information-réponse
  - Information-réponse après les éléments repris de la question

# Analyse du corpus : Synthèse

- 1 Est-il judicieux de réutiliser des éléments de la question ?
- 2 Quelle information faut-il répondre ?
- 3 Comment construire une réponse minimale ?
- 4 **La modalité (oral/écrit) influe-t-elle sur la formulation de réponse à générer ?**
- 5 Euh...et si on hésite ?
- 6 Que répondre quand on n'a pas de réponse ?
- 7 Comment répondre plusieurs informations-réponse ?

- Des réponses en général semblables
- À l'écrit : réponses plus longues et reprises par pronom plus fréquentes
- À l'oral : peu de réponses complétives et moins de réponses contenant plusieurs informations-réponse

# Analyse du corpus : Synthèse

- 1 Est-il judicieux de réutiliser des éléments de la question ?
- 2 Quelle information faut-il répondre ?
- 3 Comment construire une réponse minimale ?
- 4 La modalité (oral/écrit) influe-t-elle sur la formulation de réponse à générer ?
- 5 **Euh...et si on hésite ?**
- 6 Que répondre quand on n'a pas de réponse ?
- 7 Comment répondre plusieurs informations-réponse ?

- Des hésitations avant l'information-réponse dans 73% des cas
- Sont un indice de la confiance du locuteur

# Analyse du corpus : Synthèse

- 1 Est-il judicieux de réutiliser des éléments de la question ?
- 2 Quelle information faut-il répondre ?
- 3 Comment construire une réponse minimale ?
- 4 La modalité (oral/écrit) influe-t-elle sur la formulation de réponse à générer ?
- 5 Euh...et si on hésite ?
- 6 **Que répondre quand on n'a pas de réponse ?**
- 7 Comment répondre plusieurs informations-réponse ?

- Utilisation d'expressions figées

je ne sais pas, je n'en ai aucune idée,...

## Analyse du corpus : Synthèse

- ① Est-il judicieux de réutiliser des éléments de la question ?
- ② Quelle information faut-il répondre ?
- ③ Comment construire une réponse minimale ?
- ④ La modalité (oral/écrit) influe-t-elle sur la formulation de réponse à générer ?
- ⑤ Euh...et si on hésite ?
- ⑥ Que répondre quand on n'a pas de réponse ?
- ⑦ **Comment répondre plusieurs informations-réponse ?**

- Variation de la granularité de la réponse
- Variation du cadre de validité de la réponse

**A2525'** : au Louvre, à Paris, en France

**A60** : 28 jours les années bissextiles, 29 jours les autres

## Conclusion 1/2

### Apports

#### Construction d'un corpus de 707 questions

- dont la forme linguistique est contrôlée
- s'appuyant sur un modèle linguistique testé sur des données réelles

#### Collecte d'un corpus de 3132 réponses à des questions

- un protocole original
- un corpus unique

#### Analyse de ce corpus

- pour les systèmes de question-réponse

#### Corpus disponibles :

- 😊 bruts et annotés
- 😞 versions audio
- {*annegf, rosset, vilnat*}@limsi.fr

## Conclusion 2/2

### Limites

#### Concernant la collecte

- recrutement trop *lourd*
- 57 % des réponses du corpus sont succinctes

#### Concernant les questions

- seulement 3 types sémantiques

#### Concernant les analyses

- patrons de réponse
- mesure de la durée

# Perspectives

Implémentation au sein d'un système

Systèmes disponibles au LIMSI : RITEL, FIDJI, QAVAL,...

Évaluation de réponses en langue naturelle

- Évaluation en interaction ou en “différé” ?
- Quelle serait une tâche pour une campagne d'évaluation ?
- Quelles mesures ?

Merci !

Questions ?

# Plan

- 1 Introduction
- 2 Protocole expérimental
- 3 Construction du corpus de questions
- 4 Collecte d'un corpus de réponses
- 5 Analyses du corpus
- 6 Conclusion et perspectives