

# Generating Intelligent Numerical Answers in a Question-Answering System

Véronique Moriceau

Institut de Recherche en Informatique de Toulouse  
118, route de Narbonne, 31062 Toulouse, France  
moriceau@irit.fr

## Abstract

In this paper, we present a **question-answering** system on the Web which aims at generating intelligent answers to numerical questions. These answers are generated in a cooperative way: besides a direct answer, comments are generated to explain to the user the variation of numerical data extracted from the Web. We present the content determination and realisation tasks. We also present some elements of evaluation with respect to end-users.

## 1 Introduction

Search engines on the Web and most existing question-answering (QA) systems provide the user with a set of hyperlinks and/or Web page extracts containing answer(s) to a question. These answers may be incoherent to a certain degree: they may be equivalent, complementary, contradictory, at different levels of precision or specificity, etc. It is then quite difficult for the user to know which answer is the correct one. Thus, an analysis of relevance and coherence of candidate answers is essential.

### 1.1 Related work

Search engines on the Web produce a set of answers to a question in the form of hyperlinks or page extracts, ranked according to content or popularity criteria (Salton, 1989; Page et al., 1998). Some QA systems on the Web use other techniques: candidate answers are ranked according to a score which takes into account lexical relations between questions and answers, semantic categories of concepts, distance between words, etc. (Moldovan et al., 2003), (Narayanan and Harabagiu, 2004), (Radev and McKeown, 1998).

Recently, advanced QA systems defined relationships (equivalence, contradiction, ...) between Web page extracts or texts containing possible answers in order to combine them and to produce a single answer (Radev and McKeown, 1998), (Harabagiu and Lacatusu, 2004), (Webber et al., 2002).

Most systems provide the user with either a set of potential answers (ranked or not), or the "best" answer according to some relevance criteria. They do not provide answers which take into account **information from a set of candidate answers** or answer inconsistencies. As for logical approaches used for database query, they are based on majority approach or on source reliability. But, contrary to the assumption of (Motro et al., 2004), we noted that reliability information (information about the author, date of Web pages, ...) is rather difficult to obtain, so we assume that all Web pages are equally reliable.

### 1.2 Motivations and goals

Our framework is advanced QA systems over open domains. Our main goals are to model and to evaluate a system which, from a factoid question in natural language (**in French**), selects a set of candidate answers on the Web and generates cooperative answers in natural language. Our challenge is (1) to generate a synthetic answer instead of a list of potential answers (in order to avoid providing the user with too much information), and (2) to generate relevant comments which explain the variety of answers extracted from the Web (in order to avoid misleading the user) (Grice, 1975). In a cooperative perspective, we propose an approach for answer generation which uses answer **integration**. When several possible answers are extracted from the Web, the goal is to define a coherent core

from candidate answers and to generate a **cooperative answer**, i.e. an answer with explanations.

In this paper, we focus on the integration of numerical data in order to generate natural language cooperative answers to numerical questions. We first present some motivational problems for the generation of numerical answers in a QA system. Then, we present the content determination and realization processes. Finally, we give some elements of evaluation of our system outputs, with respect to end-users.

## 2 On numerical data

We focus on the integration of numerical data for the generation of natural language cooperative numerical answers. We first present some related work on generation from numerical data sets. Then we propose a model for the generation of cooperative numerical answers.

### 2.1 Related work

The generation of summaries from numerical data has been developed in some NLG systems. For example, the system ANA (Kukich, 1983) generates stock market reports by computing fluctuations for a day. FoG (Goldberg et al, 1994) produces weather forecasts from forecast data. More recently, StockReporter (Dale, 2003) was developed to generate summaries describing how a stock performs over a period. Yu et al. (2005) propose a system which generates summaries of sensor data from gas turbines.

Those systems have input data analysis components which are more or less efficient and describe numerical time-series data. In the framework of QA systems, there are other major problems that the previous systems do not deal with. When a numerical question is submitted to a QA system, a set of numerical data is extracted from the Web. Then, the goal is not to describe the whole data set but to find an appropriate answer, dealing with the user expectations (for example, constraints in the question) or data inconsistencies. Another important point is the analysis of numerical input data in order to identify causes (besides time) of variation.

### 2.2 A typology of numerical answers

Our challenge is to develop a formal framework for the integration of numerical data extracted from Web pages in order to produce cooperative numerical answers.

To define the different types of numerical answers, we collected a set of 80 question-answer pairs about prices, quantities, age, time, weight, temperature, speed and distance. The goal is to identify for each question-answer pair why extracted numerical values are different (is this an inconsistency? an evolution?).

A numerical question may accept several answers when numerical values vary according to some criteria. Let us consider the following examples.

#### Example 1 :

*How many inhabitants are there in France?*

- *Population census in France (1999): 60184186.*
- *61.7: number of inhabitants in France in 2004.*

#### Example 2 :

*What is the average age of marriage of women in 2004?*

- *In Iran, the average age of marriage of women was 21 years in 2004.*
- *In 2004, Moroccan women get married at the age of 27.*

#### Example 3 :

*At what temperature should I serve wine?*

- *Red wine must be served at room temperature.*
- *Champagne: between 8 and 10 °C.*
- *White wine: between 8 and 11 °C.*

The corpus analysis allows us to identify 3 main variation criteria, namely *time* (ex.1), *place* (ex.2) and *restriction* (ex.3: restriction on the focus, for example: Champagne/wine). These criteria can be combined: some numerical values vary according to time and place, to time and restrictions, etc. (for example, the average age of marriage vary according to time, place and restrictions on men/women).

### 2.3 A model for cooperative numerical answer generation

The system has to generate an answer **from a set of numerical data**. In order to identify the different problems, let us consider the following example :

*What is the average age of marriage in France?*

- *In 1972, the average age of marriage was 24.5 for men and 22.4 for women. In 2005, it is 30 for men and 28 for women.*
- *The average age of marriage in France increased from 24.5 to 26.9 for women and from 26.5 to 29 for men between 1986 and 1995.*

This set of potential answers may seem incoherent but their internal coherence can be made apparent once a variation criterion is identified. In a cooperative perspective, an answer can be for example:

*In 2005, the average age of marriage in France was 30 for men and 28 for women.*

*It increased by about 5.5 years between 1972 and 2005.*

This answer is composed of:

1. a direct answer to the question,
2. an explanation characterizing the variation mode of the numerical value.

To generate this kind of answer, it is necessary (1) to integrate candidate answers in order to elaborate a direct answer (for example by solving inconsistencies), and (2) to integrate candidate answers characteristics in order to generate an explanation.

Figure 1 presents the general architecture of our system which allows us to generate answers and explanations from several different numerical answers. Questions are submitted in natural language to QRISTAL<sup>1</sup> which analyses them and selects potential answers from the Web. Then, a grammar is applied to extract information needed for the generation of an appropriate cooperative answer. This information is mainly:

- the searched numerical value (*val*),
- the *unit* of measure,
- the question *focus*,
- the *date* and *place* of the information,
- the *restriction(s)* on the question focus ,
- the *precision* of the numerical value (for example adverbs or prepositions such as in *about 700*, ...),
- linguistic clues indicating a *variation* of the value (temporal adverbs, verbs of change/movement as in *the price increased to 200 euro*).

For the extraction of restrictions, a set of basic properties is defined (colors, form, material, etc.). Ontologies are also necessary. For example, for the question *how many inhabitants are there in France?*, population of *overseas regions* and *metropolitan* population are restrictions of *France* because they are daughters of the concept *France* in the ontology. On the contrary, prison population of France is not a restriction because *prison* is not a daughter of *France*. Several ontologies are available<sup>2</sup> but the lack of available knowledge for

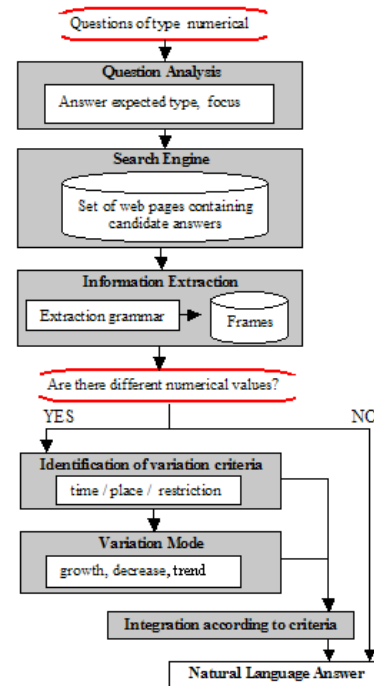


Figure 1: Architecture

some domains obviously influences the quality of answers.

We define the set  $A = \{a_1, \dots, a_N\}$ , with  $a_i$  a frame which gathers all this information for a numerical value. Figure 2 shows an extraction result.

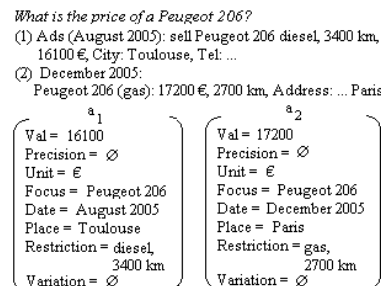


Figure 2: Extraction results

From the frame set, the variation criteria and mode of the searched numerical value are identified: these components perform content determination. Finally, a natural language answer is generated explaining those characteristics. Each of these stages is presented in the next sections.

### 3 Content determination for explanations

In order to produce explanations for data variation, the system must have a data analysis component

<sup>1</sup>www.qristal.fr, Synapse Développement

<sup>2</sup>http://www.daml.org/ontologies/

which can infer, from extracted information, the variation phenomena, criteria and mode.

### 3.1 Variation criteria

Once we have the frames representing the different numerical values, the goal is to determine if there is a variation and to identify the variation criteria of the value. We assume that there is a variation if there is at least  $k$  different numerical values with different criteria (time, place, restriction) among the  $N$  frames (for the moment, we arbitrarily set  $k = N/4$ , but this has to be evaluated). Thus, a numerical value varies according to:

1. **time** if  $T = \{a_i(Val), \exists a_j \in A,$   
such as  $a_i(Val) \neq a_j(Val)$   
 $\wedge a_i(Unit) = a_j(Unit)$   
 $\wedge a_i(Date) \neq a_j(Date) \}$   
 $\wedge card(T) \geq k$
2. **place** if  $P = \{a_i(Val), \exists a_j \in A,$   
such as  $a_i(Val) \neq a_j(Val)$   
 $\wedge a_i(Unit) = a_j(Unit)$   
 $\wedge a_i(Place) \neq a_j(Place) \}$   
 $\wedge card(P) \geq k$
3. **restriction** if  $Rt = \{a_i(Val), \exists a_j \in A,$   
such as  $a_i(Val) \neq a_j(Val)$   
 $\wedge a_i(Unit) = a_j(Unit)$   
 $\wedge a_i(Restriction) \neq a_j(Restriction) \}$   
 $\wedge card(Rt) \geq k$
4. **time and place** if  $(1) \wedge (2)$
5. **time and restriction** if  $(1) \wedge (3)$
6. **place and restriction** if  $(2) \wedge (3)$
7. **time, place and restriction** if  $(1) \wedge (2) \wedge (3)$

Numerical values can be compared only if they have the same unit of measure. If not, they have to be converted. More details about comparison rules are presented in (Moriceau, 2006).

### 3.2 Variation mode

In the case of numerical values varying over time, it is possible to characterize more precisely the variation. The idea is to draw a trend (increase, decrease, ...) of variation over time so that a precise explanation can be generated. For this purpose, we draw a regression line which determines the relationship between the two extracted variables *value* and *date*.

In particular, Pearson's correlation coefficient ( $r$ ),

related to the line slope, reflects the degree of linear relationship between two variables. It ranges from  $+1$  to  $-1$ . For example, figure 3 shows that a positive Pearson's correlation implies a general increase of values whereas a negative Pearson's correlation implies a general decrease. On the contrary, if  $r$  is low ( $-0.6 < r < 0.6$ ), then we consider that the variation is random (Fisher, 1925).

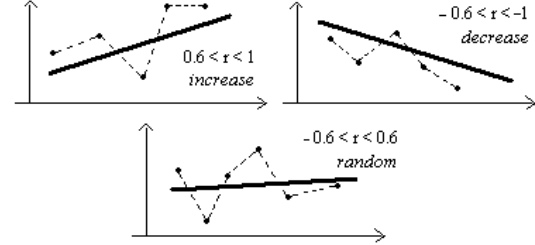


Figure 3: Variation mode

Figure 4 shows the results for the question *How many inhabitants are there in France?* Different numerical values and associated dates are extracted from Web pages. The Pearson's correlation is 0.694 meaning that the number of inhabitants increases over time (between 1999 and 2005).

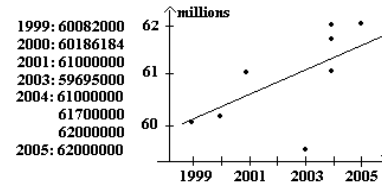


Figure 4: Variation mode: *How many inhabitants are there in France?*

## 4 Answer generation

Once the searched numerical values have been extracted and characterized by their variation criteria and mode, a cooperative answer is generated in natural language. It is composed of two parts:

- a direct answer if available,
- an explanation of the value variation.

### 4.1 Direct answer generation

#### 4.1.1 Question constraints

The content determination process for the direct answer generation is mainly guided by constraints which may be explicit or implicit in the question. For example, in the question *how many inhabitants are there in France in 2006?*, there

are explicit constraints on time and place. On the contrary, in *how many inhabitants are there in France?*, there is no constraint on time. Let  $C$  be the set of question constraints:  $C = \{C_t, C_p, C_r\}$  with :

- $C_t$ : constraint on time ( $C_t \in \{exp\_time, \emptyset\}$ ),
- $C_p$ : constraint on place ( $C_p \in \{exp\_place, \emptyset\}$ ),
- $C_r$ : constraint on restrictions ( $C_r \in \{exp\_restr, \emptyset\}$ ).

For example, in the question *what is the average age of marriage in France?*:  $C_t = \emptyset$ ,  $C_p = \text{France}$  and  $C_r = \emptyset$ .

When there is no explicit constraint in the question, we distinguish several cases:

- if there is no explicit constraint on time in the question and if a numerical variation over time has been inferred from the data set, then we assume that the user wants to have the most recent information:  $C_t = \max(\{a_i(date), a_i \in A\})$ ,
- if there is no explicit constraint on place in the question and if a numerical variation according to place has been inferred from the data set, then we assume that the user wants to have the information for the closest place to him (the system can have this information for example via a user model),
- if there is no explicit constraint on restrictions in the question and if a numerical variation according to restrictions has been inferred from the data set, then we assume that the user wants to have the information for any restrictions.

For example, on figure 5:  $C_t = 2000$  (the most recent information),  $C_p = \text{France}$  and  $C_r = \emptyset$ .

#### 4.1.2 Candidate answers

Candidate frames for direct answers are those which satisfy the set of constraints  $C$ . Let  $AC$  be the set of frames which satisfy  $C$  (via subsumption):

$$AC = \{a_i \in A, \text{ such as } \\ a_i(date) = (C_t \vee \emptyset) \wedge a_i(place) = (C_p \vee \emptyset) \wedge \\ a_i(restriction) = \begin{cases} C_r \vee \emptyset & \text{if } C_r \neq \emptyset \\ exp\_rest \vee \emptyset & \text{if } C_r = \emptyset \end{cases}$$

For figure 5:  $AC = \{a_1, a_2, a_3, a_4, a_5, a_6\}$ .

#### 4.1.3 Choosing a direct answer

A direct answer has to be generated from the set  $AC$ . We define subsets of  $AC$  which contain frames having the same restrictions: a direct answer will be generated for each relevant restriction. Let  $\mathcal{A}$  be the subsets of frames satisfying the question constraints and having the same restrictions:  $\mathcal{A} = \{AC_1, \dots, AC_M\}$  with:

$$AC_i = \{a_j, \text{ such as } \forall a_j, a_k \in AC, \\ a_j(restriction) = a_k(restriction) \\ \vee a_j(restriction) = \emptyset\}, \\ \text{and } AC_1, \dots, AC_M \text{ are disjoint.}$$

For figure 5:  $\mathcal{A} = \{AC_1, AC_2\}$  with:

$AC_1 = \{a_1, a_3, a_5\}$ , subset for restriction *women*,  
 $AC_2 = \{a_2, a_4, a_6\}$ , subset for restriction *men*.

Then, for each element in  $\mathcal{A}$ , an answer is generated :

$$\forall AC_i \in \mathcal{A}, \text{ answer} = \text{generate\_answer}(AC_i).$$

Each element of  $\mathcal{A}$  may contain one or several frames, i.e. one or several numerical data. Some of these values may be aberrant (for example, *How high is the Eiffel Tower?* 300m, 324m, 18cm): they are filtered out via classical statistical methods (use of the standard deviation). Among the remaining frames, values may be equal or not at different degrees (rounded values, for example). Those values have to be integrated so that a synthetic answer can be generated.

There are many operators used in logical approaches for fusion: conjunction, disjunction, average, etc. But, they may produce an answer which is not cooperative: a conjunction or disjunction of all candidates may mislead users; the average of candidates is an "artificial" answer since it has been computed and not extracted from Web pages.

Our approach allows the system to choose a value among the set of possible values, dealing with the problem of rounded or approximative data. Candidate values are represented by an oriented graph whose arcs are weighted with the cost between the two linked values and the weight ( $w$ ) of the departure value (its number of occurrences). A graph  $\mathcal{G}$  of numerical values is defined by  $\mathcal{N}$  the set of nodes (set of values) and  $\mathcal{Arc}$  the set of arcs. The cost  $c(x, y)$  of  $arc(x, y)$  is:

$$\frac{|x - y|}{y} \times (w(x) + \sum_{i=1}^n w(x_i)) + \sum_{i=1}^n c(x_i, x).$$

with  $(x_1, \dots, x_n, x)$  a path from  $x_1$  to  $x$ .

Finally, we define a fusion operator which selects the value which is used for the direct answer. This value is the one which maximizes the difference ( $cost(x)$ ) between the cost to leave this value and the cost to arrive to this value:

$\begin{pmatrix} a_1 \\ \text{Val} = 27.7 \\ \text{Precision} = \emptyset \\ \text{Unit} = \text{years} \\ \text{Focus} = \text{age of marriage} \\ \text{Date} = 01/01/2000 \\ \text{Place} = \text{France} \\ \text{Restriction} = \text{women} \\ \text{Variation} = \emptyset \end{pmatrix}$	$\begin{pmatrix} a_2 \\ \text{Val} = 29.8 \\ \text{Precision} = \emptyset \\ \text{Unit} = \text{years} \\ \text{Focus} = \text{age of marriage} \\ \text{Date} = 01/01/2000 \\ \text{Place} = \text{France} \\ \text{Restriction} = \text{men} \\ \text{Variation} = \emptyset \end{pmatrix}$	$\begin{pmatrix} a_3 \\ \text{Val} = 28 \\ \text{Precision} = \emptyset \\ \text{Unit} = \text{years} \\ \text{Focus} = \text{age of marriage} \\ \text{Date} = 01/01/2000 \\ \text{Place} = \emptyset \\ \text{Restriction} = \text{women} \\ \text{Variation} = \emptyset \end{pmatrix}$	$\begin{pmatrix} a_4 \\ \text{Val} = 30 \\ \text{Precision} = \emptyset \\ \text{Unit} = \text{years} \\ \text{Focus} = \text{age of marriage} \\ \text{Date} = 01/01/2000 \\ \text{Place} = \emptyset \\ \text{Restriction} = \text{men} \\ \text{Variation} = \emptyset \end{pmatrix}$	$\begin{pmatrix} a_5 \\ \text{Val} = 28.5 \\ \text{Precision} = \emptyset \\ \text{Unit} = \text{years} \\ \text{Focus} = \text{age of marriage} \\ \text{Date} = \emptyset \\ \text{Place} = \emptyset \\ \text{Restriction} = \text{women} \\ \text{Variation} = \emptyset \end{pmatrix}$
$\begin{pmatrix} a_6 \\ \text{Val} = 30.6 \\ \text{Precision} = \emptyset \\ \text{Unit} = \text{years} \\ \text{Focus} = \text{age of marriage} \\ \text{Date} = \emptyset \\ \text{Place} = \emptyset \\ \text{Restriction} = \text{men} \\ \text{Variation} = \emptyset \end{pmatrix}$	$\begin{pmatrix} a_7 \\ \text{Val} = 25.8 \\ \text{Precision} = \emptyset \\ \text{Unit} = \text{years} \\ \text{Focus} = \text{age of marriage} \\ \text{Date} = 1990 \\ \text{Place} = \emptyset \\ \text{Restriction} = \text{women} \\ \text{Variation} = \emptyset \end{pmatrix}$	$\begin{pmatrix} a_8 \\ \text{Val} = 27.8 \\ \text{Precision} = \emptyset \\ \text{Unit} = \text{years} \\ \text{Focus} = \text{age of marriage} \\ \text{Date} = 1990 \\ \text{Place} = \emptyset \\ \text{Restriction} = \text{men} \\ \text{Variation} = \emptyset \end{pmatrix}$	$\begin{pmatrix} a_9 \\ \text{Val} = 24.2 \\ \text{Precision} = \emptyset \\ \text{Unit} = \text{years} \\ \text{Focus} = \text{age of marriage} \\ \text{Date} = 1985 \\ \text{Place} = \text{France} \\ \text{Restriction} = \text{women} \\ \text{Variation} = \emptyset \end{pmatrix}$	$\begin{pmatrix} a_{10} \\ \text{Val} = 26.3 \\ \text{Precision} = \emptyset \\ \text{Unit} = \text{years} \\ \text{Focus} = \text{age of marriage} \\ \text{Date} = 1985 \\ \text{Place} = \text{France} \\ \text{Restriction} = \text{men} \\ \text{Variation} = \emptyset \end{pmatrix}$

Figure 5: Data set for *What is the average age of marriage in France?*

answer =  $y \in \mathcal{N}$ , such as

$$\text{cost}(y) = \max(\{ \text{cost}(n), \forall n \in \mathcal{N}, \\ \text{cost}(n) = \text{cost\_leave}(n) - \text{cost\_arrive}(n) \})$$

with:  $\text{cost\_leave}(x) = \sum_i c(x, x_i)$  and,  
 $\text{cost\_arrive}(x) = \sum_i c(x_i, x)$ .

Let us consider an example. The following values are candidate for the direct answer to the question *How high is the Mont-Blanc?*: 4800, 4807 (2 occurrences), 4808 (2 occurrences), 4808.75, 4810 (8 occurrences) and 4813. Figure 6 shows the graph of values: in this example, the value which maximizes the costs is 4810.

From the selected value, the system generates a direct answer in natural language in the form of *Focus Verb (Precision) Value*. For example, the generated answer for *How high is the Mont-Blanc?* is *The Mont-Blanc is about 4810 meters high*. Here the preposition *about* indicates to the user that the given value is an approximation. For the question *what is the average age of marriage in France?*, a direct answer has to be generated for each restriction. For the restriction *men* ( $AC_2$ ), there are 3 candidate values: 29.8, 30 and 30.6, the value which minimizes the costs being 30. For the restriction *women* ( $AC_1$ ), there are also 3 candidate values: 27.7, 28 and 28.5, the value which minimizes the costs being 28. After aggregation process, the generated direct answer is: *In 2000, the average age of marriage in France was about 30 years for men and 28 years for women*.

## 4.2 Explanation generation

The generation of the cooperative part of the answer is complex because it requires lexical knowledge. This part of the answer has to explain to the user variation phenomena of search values: when a variation of values is identified and char-

acterised, an explanation is generated in the form of *X varies according to Criteria*. In the case of variation according to restrictions or properties of the focus, a generalizer is generated. For example, the average age of marriage varies for men and women: the explanation is in the form *the average age of marriage varies according to sex*. The generalizer is the mother concept in the ontology or a property of the mother concept (Benamara, 2004). For numerical value varying over time, if the variation mode (increase or decrease) is identified, a more precise explanation is generated: *X increased/decreased between... and...* instead of *X varies over time*.

Here, verbs are used to express precisely numerical variations. The lexicalisation process needs deep lexical descriptions. We use for that purpose a classification of French verbs (Saint-Dizier, 1999) based on the main classes defined by WordNet. The classes we are interested in for our task are mainly those of verbs of state (*have, be, weight*, etc.), verbs of change (*increase, decrease*, etc.) and verbs of movement (*climb, move forward/backward*, etc.) used metaphorically (Moriceau et al, 2003). From these classes, we selected a set of about 100 verbs which can be applied to numerical values.

From these classes, we characterized sub-classes of growth, decrease, etc., so that the lexicalisation task is constrained by the type of verbs which has to be used according to the variation mode.

A deep semantics of verbs is necessary to generate an answer which takes into account the characteristics of numerical variation as well as possible: for example, the variation mode but also the speed and range of the variation. Thus, for each sub-class of verbs and its associated variation mode, we need a refined description of ontological domains and selectional restrictions so that



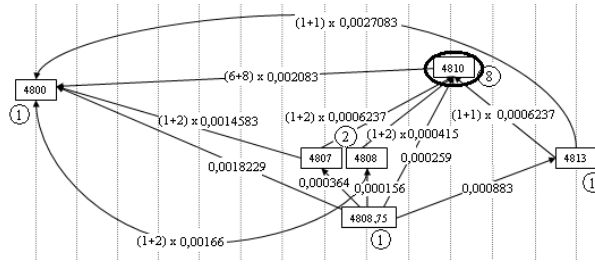


Figure 6: Graph of candidate values for *How high is the Mont-Blanc?*

an appropriate verb lexicalisation can be chosen: which verb can be applied to prices, to age, etc.? (Moriceau et al, 2003). We propose to use proportional series representing verb sub-classes according to the speed and amplitude of variation. For example, the use of *climb* (resp. *drop*) indicates a faster growth (resp. decrease) than *go up* (resp. *go down*): the verb *climb* is preferred for the generation of *Gas prices climb 20.3% in october 2005* whereas *go up* is preferred in *Gas prices went up 7.2% in september 2005*.

Verbs can possibly be associated with a preposition that refines the information (*The average age of marriage increased by about 5.5 years between 1972 and 2005*).

### 4.3 Answer justification

Our system generates a cooperative answer composed of a direct answer to the question and an explanation for the possible variation of the searched numerical value. But the answer may not be sure because of a too high/low number of candidate values to the direct answer. In this case, it may be useful to add some additional information for the user in order to justify or complete the generated answer.

We propose to add a **know-how** component to our system, which provides the user with one or two relevant Web page extracts besides the generated answer whenever it is necessary. These extracts must contain information about the searched numerical values, and for example some explanations of the causes of numerical variation. Some linguistic clues can be used to select page extracts: number of numerical values concerning the question focus, causal marks (because of, due to, ...), etc. Figure 7 shows an output example of our system.

Question: **How high is the Mont-Blanc?**

Answer: The Mont-Blanc is about 4810 metres high.  
The height of the Mont-Blanc varies over time.

Some more details:

Mont Blanc in the Alps, is the highest mountain in Western Europe. Its height is about 4810 m, but varies from year to year by a few metres, depending on snowfall and climate conditions.  
Mont Blanc has traditionally been considered to be 4807 m high, but GPS-based measurements made in 2002 by the Institut Géographique National and other experts were 4810,40 m. After the heat wave in 2003, new measurements made on September, 6<sup>th</sup> and 7<sup>th</sup> were 4808,45 m. These seem to result from fluctuations, caused by the weather. This interpretation is disputed, since the 2003 heat wave didn't significantly affect the glaciers above 4000 meters altitude. It could also be explained by random movements of the summit icecap, due to the violent winds at this altitude.

Figure 7: An output example

## 5 Evaluation

In this section, we present some elements of evaluation of our system with respect to 15 end-users<sup>3</sup>.

We first evaluated how users behave when they are faced with different candidate answers to a question. To each user, we presented 5 numerical questions and their candidate answers which vary according to time or restrictions and ask them to produce their own answer from candidate answers. For numerical answers varying according to restrictions, 93% of subjects produce answers explaining the different numerical values for each restriction. For numerical answers varying over time, 80% of subjects produce answers giving the most recent information (20% of subjects produce an answer which a summary of all candidate values). This validates our hypothesis presented in section 4.1.1.

The second point we evaluated is the answer order. Our system produces answers in the form of a direct answer, then an explanation and a justification (page extract) if necessary. We proposed to users answers with these three parts arranged randomly. Contrary to (Yu et al, 2005) which propose first an overview and then a zoom on inter-

<sup>3</sup>Subjects are between 20 and 35 years old and are accustomed to using search engines.

esting phenomena, 73% of subjects preferred the order proposed by our system, perhaps because, in QA systems, users want to have a direct answer to their question before having explanations.

The last point we evaluated is the quality of the system answers. For this purpose, we asked subjects to choose, for 5 questions, which answer they prefer among: the system answer, an average, an interval and a disjunction of all candidate answers. 91% of subjects preferred the system answer. 75% of subjects found that the explanation produced is useful and only 31% of subjects consulted the Web page extract (28% of these found it useful).

## 6 Conclusion

We proposed a question-answering system which generates intelligent answers to numerical questions. Candidate answers are first extracted from the Web. Generated answers are composed of three parts: (1) a direct answer: the content determination process "chooses" a direct answer among candidates, dealing with data inconsistencies and approximations, (2) an explanation: the content determination process allows to identify, from data sets, the possible value variations and to infer their variation criteria (time, place or restrictions on the question focus), and (3) a possible Web page extract. This work has several future directions among which we plan:

- to define precisely in which cases it is useful to propose a Web page extract as a justification and,
- to measure the relevance of restrictions on the question focus to avoid generating an enumeration of values corresponding to irrelevant restrictions.

## References

- F. Benamara. 2004. Generating Intensional Answers in Intelligent Question Answering Systems. *LNAI Series*, volume 3123, Springer.
- R. Dale. 2003. <http://www.ics.mq.edu.au/~lgt-demo/StockReporter/>.
- R. A. Fisher 1925. *Statistical Methods for Research Workers*, originally published in London by Oliver and Boyd.
- E. Goldberg, N. Driedger, R. Kittredge. 1994. Using natural language processing to produce weather forecasts. *IEEE Expert* 9(2).
- H.P. Grice. 1975. Logic and conversation. In P. Cole and J.L. Morgan, (eds.): *Syntax and Semantics*, Vol. 3, Speech Acts, New York, Academic Press.
- S. Harabagiu and F. Lacatusu. 2004. *Strategies for Advanced Question Answering*. Proceedings of the Workshop on Pragmatics of Question Answering at HLT-NAACL 2004.
- K. Kukich. 1983. Knowledge-based report generation: a knowledge engineering approach to natural language report generation. Ph.D. Thesis, Information Science Department, University of Pittsburgh.
- D. Moldovan, C. Clark, S. Harabagiu and S. Maiorano. 2003. *COGEX: A Logic Prover for Question Answering*. Proceedings of HLT-NAACL 2003.
- V. Moriceau and P. Saint-Dizier. 2003. *A Conceptual Treatment of Metaphors for NLP*. Proceedings of ICON, Mysore, India.
- V. Moriceau. 2006. *Numerical Data Integration for Question-Answering*. Proceedings of EACL-KRAQ'06, Trento, Italy.
- A. Motro, P. Anokhin. 2004. Fusionplex: resolution of data inconsistencies in the integration of heterogeneous information sources. *Information Fusion*, Elsevier.
- S. Narayanan, S. Harabagiu. 2004. *Answering Questions Using Advanced Semantics and Probabilistic Inference*. Proceedings of the Workshop on Pragmatics of Question Answering, HLT-NAACL, Boston, USA, 2004.
- L. Page, S. Brin, R. Motwani, T. Winograd. 1998. *The PageRank Citation Ranking: Bringing Order to the Web*. Technical Report, Computer Science Department, Stanford University.
- D.R. Radev and K.R. McKeown. 1998. Generating Natural Language Summaries from Multiple On-Line Sources. *Computational Linguistics*, vol. 24, issue 3 - Natural Language Generation.
- P. Saint-Dizier. 1999. Alternations and Verb Semantic Classes for French. *Predicative Forms for NL and LKB*, Kluwer Academic.
- P. Saint-Dizier. 2005. *PrepNet: a Framework for Describing Prepositions: preliminary investigation results*. Proceedings of IWCS'05, Tilburg, The Netherlands.
- G. Salton. 2002. Automatic Text Processing. *The Transformation, Analysis and Retrieval of Information by Computer*, Addison-Wesley.
- B. Webber, C. Gardent and J. Bos. 2002. *Position statement: Inference in Question Answering*. Proceedings of LREC, Las Palmas, Spain.
- J. Yu, E. Reiter, J. Hunter, C. Mellish. 2005. Choosing the content of textual summaries of large time-series data sets. *Natural Language Engineering*, 11.