Internal Use Only

TR-I-0044

Record of Six Work Sessions on Concepts, Methods, and Tools from Existing Running Real-Size MT Systems

> ワークセッション記録: MTシステムの概念、手法、ツール

> > Christian Boitet クリスチャン・ブワテ

> > > October, 1988

ATR Interpreting Telephony Research Laboratories ATR 自動翻訳電話研究所

### Foreword

This is a record of six work sessions organized by Professor Christian Boitet during his stay at ATR.

Using his long research experience in the field of Machine Translation, Professor Boitet reviewed the concepts, methods, and tools of existing or previously running real-size MT systems. He discussed them with the ATR researchers in order to work out the "best set" of principles for an Automatic Interpretation System, taking into account existing dangers, illusions, and solutions.

He prepared more than a hundred transparent sheets for the work sessions, which are the main sources of this record. However, some errors might have occured in the editing process of the hand-written transparent sheets into this typed format. I would like to apologize to the author for the remaining errors.

(T. Aizawa)

# Six Work Sessions on Concepts, Methods, and Tools from Existing Running Real-Size MT Systems

### **Christian Boitet**

ATR Interpreting Telephony Research Laboratories, Osaka and GETA (Study Group for Machine Translation) Joseph Fourier University & CNRS, Grenoble

# TABLE OF CONTENTS

	INTRODUCTION	2
	Why refer to MT systems if we attack MI systems?	
1.	Historical Perspective	3
	Through a retracing of Prof. B. Vauquois' contribution	
2.	Basic Architectures and Concepts	13
	A first level of description should also apply to MI systems	
3.	Representation of the Units of Translation	21
	Input / output objects	
4.	Specification & Application of Grammars & Dictionaries	37
	How to build linguistic processors	
5.	Specialized Languages for Linguistic Programming	43
	The tools of the trade	
6.	Organization of the Whole MT / MI Ststem	56
	Unified architecture rather than accidental pasting	

### INTRODUCTION

Why refer to MT (Machine Translation) systems if we attack MI (Machine Interpretation) systems?

- Prof. Nagao's warning not to follow MT design blindly
- Written & spoken language, translation & interpetation: A common basis
- Research in NLP as a finite-state device much has been done .....
   .....and forgotten
- Usual scientific practice to build on accumulated knowledge .....
   .....without rediscovering the wheel
- B. Vauquois' contribution & evolution;
  - Coincidental with the history of modern MT
  - Illustrates the bychological side

1st Work Session (April 21, 1988)

# **Historical Perspective**

# **Outline of Presentations**

### A. Through a retracing of B. Vauquois' contribution & evolution

1. The researcher & ms rucas	1.	The research	ıer & l	hisideas
------------------------------	----	--------------	---------	----------

1960~65	Applying formal language theory
1965~70	Integrating modern semantic theories
1970~75	Defining modular grammars and heuristic methods
1975~80	Inventing multilevel (M-)Structures
1980~85	"Static grammars" for specifying the string-tree
	correspondences
1985~	Towards a new treatment of ambiguity

2. The implementer & his methods

Erudition, eclectism & pragmatism in linguistics specialized languages adapted to the tasks emergence of lingware engineering methods

- B. With a rapid survey of some MT systems
  - 1. SYSTRAN & variants
  - 2. CETA & METAL
  - 3. SUSY & LOGOS
  - 4. METEO & TAUM-AVIATION
  - 5. ARIANE-78 & -85
  - 6. Japanese systems
  - 7. EUROTRA?
- C. A summary of the important issues

### 1960~65 Applying formal language theory

 $\begin{array}{c} \text{Grammar in CNF} \\ \text{A} \rightarrow \text{BC} \\ \text{A} \rightarrow \text{t} \end{array}$ 

١

With attributes on symbols

- Normal (gender, number .....)

- "Vehicular" (for discontinuines)

And validations/saturations on rules

 $\begin{array}{rll} \mathrm{R}_{A1} & : & A \rightarrow \mathrm{BC} & / \mathrm{P}(\alpha(\mathrm{B}), \alpha(\mathrm{C})) \\ & & / \alpha(\mathrm{A}) & : = \mathrm{F}(\alpha(\mathrm{B}), \alpha(\mathrm{C})) \\ & & / \mathrm{-}(\mathrm{R}_i, \mathrm{R}_j, \ldots .), + (\mathrm{R}_r, \ldots .). \end{array}$ 

Cocke algorithm



with some improvements.

### 1965~70 Integrating modern semantic theories

• Need for a "PIVOT",

More universal than binary attributed phrase-structures

- Ideas from
  - Tesnière : Actants Mel'čuk : Meaning-text Šaumjan : Hybrid pivot

Necessity to

Select one binary structure transform it into PIVOT form

### The PIVOT

A "universal" structure (grammatical & relational symbols) with a languagedependent lexicon.



John convinced him to go using charm

Selection of initial structure and obtention of PIVOT structure

• Preference rules

(For same covered substring)

Not implemented yet as metalanguage

Transformation into an initial Dependency structure



Then, transformations expressed in a special metalanguage

#### Defining modular grammars & heuristic methods $1970 \sim 75$

- Intractability of big rule systems for human control 0 Still applicable today
- Difficulty in expressing preferences
- Belief that direct methods based on heuristics would be quicker & more adequate.

Analysis in ARIANE-78 ⇒

:

Syntax / Semantics :

Morphology

**Combinatorial & Heuristics** (ATEF) Procedural & Heuristics

(ROBRA)

<u>ATEF</u>

Computation :

ND / \* / Depth-first & Heuristics



• Output : Factorizing tree



Time flies like an arrow

### ROBRA

Tree-transformational systems

• Control graph

- Ordered grammars on nodes
- Conditions on the arcs
- ND / 1 / Depth first
  - & Context-senstive rules (on trees) Recursive calls Parallel application of rules Guarded iteration & recursion

See GRADE for comparison.

### 1975~80 Inventing multilevel (M-)Structures

Layers of interpretation levels

Syntactic categories	CAT	Words
	K	Clauses
Syntactic functions	SUBJ	
	OBJ1	
	OBJ2	
Slightly different	ATROBJ	
From dependency	CPAG	

Relations which relate only words, not words / groups with words / groups.

Logical argument places

ARGO ARG1 ARG2 TRL01 TRL12

Semantic relations

INSTR AGENT CAUSE "TOY" ACCOMP LOCAL MEASURE CONSEQ QFIER QFIED Other levels

ActualisationTime / Mode / AspectDiscourse typeQuest / Indir / Cmd .....Theme/Rheme- Not implemented yetSem. featuresConcr, Abstr, Docum,Presfied, Animate,<br/>Human, .....

+ Lexical units (Sem. derivational families)

Considering M-structures as implicit generators of other structures

С-

Classical

F - Structures SPA -

Alternate structures: Encoding Ambiguities Doubts

& During transfer Encoding

Advices or Orders To the generator

• At the beginning of generation, 1 structure encodes all paraphrases legitimate (For translation)

 $\triangle \langle$ 

Check the correct working of the device

Check that the device works correctly

Usually, but not necessarily, lexical units are not changed during generation.

1980~85 Specifying string-tree correspondences by static grammars

 Necessity of specification formalism, in view of the size of the (procedural) grammars, even if modular.

- Used in practice since 1983 (national project)
- A small example on a formal language (using a simplified syntax)



**Representation** Tree

By the following derivation

 $a_2$ 

b2



See some analogy with XYZGs ... . But in Zaharin-86, Unification --Identification.

### 1985~ Toward a new treatment of ambiguity

• In practice, "static grammars" also contain (up to 75% !) preference rules



(In procedural analysis grammar, if T1 is chosen, T2 may be encoded.)

- Desire to treat ambiguities directly at the level of the final description of the linguistic computation
  - Open area, research theme
  - Somewhat analogous to more ancient techniques, but in the spirit of today's "expert systems".

### The implementor & his methods

Erudition, eclectism, pragmatism in linguistics

- Not a blind follower of some ideology ("Mainstream" or other)
- No linguistic theory actually covers more than a fraction of the linguistic facts / phenomena.
- Good ideas may come from different sources Theoreticians and grammarians
- Theories & usually accepted concept are sometime inexact and must be corrected / improved.

# Specialized languages for linguistic programming adapted to the various tasks

• No unique computational

Framework is good for everything

Very general SLLPs tend to be very inefficient for specific tasks such

as morphology

Maybe a useful consideration also in the case of MI systems.

- Building SLLPs on subrecursive formalisms, with the possibility
  - To allow for some "liberty switches" static detection
  - To evaluate the order of complexity of the associated computations.

### The emergence of lingware engineering

- "Lingware" coined at GETA to reflect the parallelism with software engineering.
- Development & maintenance problems
   Some figures derived from CAT-NP quite in accordance with B.
   Vauquois' estimates.
  - From all points of view
     Linguistics (Content: Grammar / Lexicons)
     Computer science
     Formalisms & mathematics
     Ruilding MT systems pagagitates to fix the state

Building MT systems necessitates to fix the state of the art at some point

⇒ Success of METAL Failure of eternal prototyping / EUROTRA) 2nd Work Session (May 10, 1988)

# **Basic Architectures and Concepts**

### A rapid survey of real-size MT systems

### **SYSTRAN & Variants**

Fixed flow of control

No possibility to backtrack or process several solutions in parallel

No explicit grammar rules

No .....

But

Enormous amount of work put in dictionaries

Experience on extremely large corpuses

Incredible speed

Usable for

- information-gathering in general
- translation, in very specific settings (XEROX .....)

#### CETA & METAL

- Comparable in computational technique
  - All-path of analysis
  - Attached transformations for analysis
  - Big dictionaries

But some differences

Preference rules	/	Weights		
Hybrid PIVOT		Hybrid transfer (T + G mixed)		
Dependency structures		Constituent structures		
All-or-nothing		Fail-soft		
Implementation of the SLLPs				
Assembler	/	Lisp		

- CETA faster (80000 ipw / 1.5 Mopw)
- METAL (1981~86) 15 years later, more usable & up-to-date
- Becase ceta model was much more theory-oriented .....?  $\rightarrow$  Discussion !

### SUSY & LOGOS

- Both ~ 1.5 generation:
   SLLP for dictionaries, not for grammars
   →Difficult to adapt the grammars
- Both usable (LOGOS commercially, SUSY exper. at IAI)
- Both relying on combinatorial approach for analysis
- Both unclear linguistically & computationally Structures manipulated & produced not independently defined, but in terms of the implementation.
- Both, but LOGOS even more, give great emphasis on dictionarybuilding facilities.

### **METEO & TAUM-Aviation**

METEO

Q-systems (A. Colmerauer) 95% quality since 1983

In 1985, switched to fail-safe approach

1988: 120 w/mn on a lap top 300,000 pages translated since 1977

(~ 100 p/day) (or 25,000 w/day)

Sematics-based analysis (Categories such as "Meteorological Event")

Q-systems : General "Addition" Production system (on Q-graph)

- All rules apply on all possible occurrences; arcs used in L.H.S. are marked.
- If & when application stops, marked arcs & "loose ends" are erased; another Q-graph is output.

 $\{A^nB^nC^n \mid n \ge 1\}$  in Q-Systems

Input: -01-A + A + A + B + B + B + C + C + C - 02-

G1) Parameter-Free  

$$A + B + C = = S$$

$$A + S + B(*) + C = = S$$

$$B + B = = B + B(*)$$

$$B(*) + C = = C + B(*)$$

$$(B(*) = B$$

Counting Q-Grammar (Output =  $S(1^n)$ )

G2

A + B + C	= = S(I)	
$A + S(U^*) + B(^*) + C$	= = S(I, U*)	** U* is a list.
B + B	= = B + B(*)	
B(*) + C	= = C + B(*)	

Also reversible

(produce  $A^nB^nC^n$  from  $S(1^n)$ ).

G3) Structuring Q-Grammar (Output =  $S(A_1, B_1, C_1, S(....)))$ 

> Also reversible (produce AnBnCn from S(A, B, C, S(A, B, C, S(...., S(A, B, C))))).

### TAUM-Aviation (1976/77~1981)

- Aim at large sublanguage (Maintenance manuals for aircraft)
- Several SLLP's

Q-systems	(Generation & 1st version of morph. an.)
REZO	(ATN-like transducers of R-graphs
	binary features added)
LEXTRA	(Lexical transfer + transformational rules)

• All or nothing approach

 $(Acceptor-based) \rightarrow High quality$ 

Butfailure (1981)

Main reason:

Dictionaries (transfer) too costly because of complex transformations

Also:

ATNs bad for large grammars

(Delicate heuristics & "PATCH-UPs")

Too detailed semantic feaures

 $\rightarrow$  increase in cost

 $\rightarrow$  decrease in adaptability

### ARIANE-78 & 85

- In both, structural analysis using tree-transducer (in ROBRA)
- Pattern-matching facilitates heuristic programming, but decreases speed (~ × 40)
- Dictionary aspect not strong enough in Ariane-78:
  - $\rightarrow$  Improvements in ARIANE-85 ("Lexical expansion" optional phases)
  - $\rightarrow$  External "neutral" lexical data base + support utilities (CAT-NP)
- Some combinatorial step could be added to reduce ambiguities before structural analysis.
- However, tree-transducer gives more flexibility than attached transformations or ATN-like mechanism.

### $\{AnBnCn \mid n \ge 1\}$ in ROBRA

Input  $X (A_1, ..., A_n, B_1, ..., B_n, C_1, ..., C_n)$ Control Graph:  $) \xrightarrow{S(*A, B, C, *)} S(*, A, B, C, S, *)$ 

(Slightly simplified syntax)



Tree-transducers MU HITACHI ATNs TOSHIBA

-18-

### CF-based

# IBM JAPAN KDD(KATE) NTT (LUTE)?

"Unification-based"

Dictionaries (Large)

Direct

Neutral

Conceptual

HITACHI, KDD, ..... MU (JICST) FUJITSU, NEC

Enormous work, perhaps multiplied if pivot ( $\rightarrow$  sort of normalization across languages).

Summary of some important points in MT (of written texts)

Basic architecture & concepts

- The external characteristics for decision-makers
- Linked with intended use/evolution

Representation of the units of translation

..... And what they are (sentences, paragraphs .....)

¿ Can developers understand them ?

Specification / application of grammars & dict.

The declarative/procedural

Competence / performance Controversy

Specialized languages for linguistic programming A classification by the cost of running Comparable applications

Organization of the whole MT system

The importance of size

Various environments necessary (tests, debugging, prod.) How to achieve modularity at user-lever (lingware / texts) External characteristics seen by potential users (interested in technique)

- Basic architecture & concepts
  - Direct Transfer Pivot
  - ② Structure of grammar / dictionaries
     Explicit / Buried in codes
    - Modular / Big block

Domain Typology

Expert knowledge available?

Mixed with linguistic Knowledge. (TAUM-METEO) Separated ("coupling" 2G + KB) Integrated by compilation (CMU)

③ Treatment of ambiguities

Combinatorial / Heuristic / ..... None Batch / Interactive 3rd Work Session (May 17, 1988)

# **Representation of the Units of Translation**

- In existing MT systems
  - Computer structures
  - Linguistic content / interpretation
  - Dubious / incomplete / ambiguous analysis
- What about (unmentioned) F-structures ?
   Derivation trees vs. representation trees
   Geometry (lists) vs. Algebra (sets)
   Keep (attributed) trees with F-structures, if any (& why not ?)
- A sketch of the beginning of a tentative approach for a convenient <u>RS</u> for <u>MIS</u> Some motivations

Levels of description, contents of cells Open problems (& possibilities)

### **More Internal Characteristics**

• Representation of the units of translation

① Computer structures

Strings Lattices / charts Labelled / decorated / "featured" trees General graphs / semantic networks Logical formulae

More "AI" oriented

 $\pm$  Encoding of ambiguity  $\rightarrow$  important for MT

② Linguistic content

Levels of interpretation (CAT/K, SF, LR, SR) Grammatical properties Actualisation, determination, quantification Semantic features / restrictions Discourse-related (theme, emphasis, .....)

③ Representation of ambiguity /doubt .....
 ± Encoding of strategic indications for later processing

### What Units of Translation?

- In all systems but GETA's, sentence by sentence
  - $\rightarrow$  Quite empty talk about an aphora resolution, discourse structure, etc.
- At TAUM (~ 1972-73), some experiments with "transferring" to next .....
   → Not quite useful
- In GETA's systems, (~ 1/2 page or more), still
  - $\rightarrow$  Only used for intersentential anaphora, sentence splitting / joining
  - → Not enough linguistic knowledge for really more discourse-related proc.
  - → Long term memory needed e.g. definition of acronyms at beginning
- In MIS? One can suggest
  - Current complete utterance + All past dialogue

+ Abstractions (SPK, TOPIC)

### Here, we center on the representation of the current occurrence

L E	HYP.	In dialogue,	$\leq$	4/5 lines, or 40/50 words
N G		In conversation,	$\leq$	15/20 lines, or 150/200 words
$\left. {{ m T}\over { m H}} \right $		In meeting,	$\leq$	Maximum allowed by the chairman !

Anyway, a decomposition in paragraph-like units seems reasonable. Then use previous limit.

2	Input or output representations	
Ķ	"Interface Structures"	
I N	Internal representations	
D S	"Working Structures"	

# Elements / Constructors of Structures

Strings	of charact of comple	ters (not so simple) x elements ("Nodes")	
Binary features	+ HUM - HUM	+CONCR +ANIM +LOC +PR -CONCR -ANIM -LOC -PR	SN SN
		± ± ± City, country,	±
Typed attributes	GNR : SF : LEV :	<u>nex</u> (MAS, FEM, NEU) <u>exc</u> (SUBJ, OBJ1, OBJ2, CPAG, <u>exc</u> (RESPECT, POLITE, NEUTRA FAMILLAR, INJURIOUS)	) L,
	CAT :	$\underline{\text{exc}}$ (N : $\underline{\text{nex}}$ (CN, PN),	

A : <u>nex</u> (ADV, ADJ), .....)

Hierarchy may be implicit ( $\rightarrow$  Prop. list)

CAT :	exc	(N, A,)
SUBN :	nex	(CN, PN)
SUBA :	nex	(ADV, ADJ)

- 23 -

Trees



Problem : Algebraic information mixed with the geometric structure

Attributed	Node: 'CITY' 🛛 [ SG, CONCR, ABSTR, N, ]
(Binary features)	Label
	Features

Decorated (Typed attributes) Node: 'NP' • [NB (SG), SEM (CONCR, ABSTR), HEAD ('CITY')]

Many systems (incl. SUSY, GETA, SALAT-Heidelberg .....) Use the label to express :

• The "Main category" on non-terminals

• The "Lexical reference" on terminals

Charts / Q-Graphs

Charts (KAY, MIND system, ~1965)

Loop-free graph with linear basis



Q-graphs (Colmerauer, Q-systems, 1970)



Objects: Labelled Trees (Q-system) Attributed Trees (TAUM-AVIA) Dotted Rules (CHART PARSERI) Etc.

Other types of interface structures

Logical formulae (McCORD, LMT) Sematic network ¿(FUJITSU, ATLAS II)? Others?

Working Structures

Same, plus

String of decorated nodes (I Pointers)

SYSTRAN

GETA

Lattice of (strings of) dec. nodes

(Morph, Analysis)





# Stacks / Arrays of trees

VI)
L)

# $\boxtimes$ Interface Structure

# 🛛 Working Structure

Struc- tures	Strings			Strings Labelled Trees				C Q-Graphs		O G t r	L F o o g r		
	Simple	Decor	ations	Simple	+Binary Features	I	ecorate	d	r t	Simple	+Binary Features	n a e p r h	1 m c u a l l a
Steps			Pointers				Pointers	Weights	S			S	е
Input Text	~ALL						-			TAUM T-AVIA			
Morph. Analysis	1				-					TAUM		ARIANE	LMT ?
		SYSTRAN				ARIANE MU			CETA	METEO TAUM-AV	ATION		LMT
Struct. Analysis			SYSTRAN			(STACK OF ARIANE MU	)		CETA METAL ~ALL CF-based	METEO	T-AVIA		LMT
Interface Source Structure			SYSTRAN (?)	METEO (1 ARC)	T-AVIA (1 ARC)	ARIANE MU METAL	CETA			METEO	T-AVIA	FUJITSU (?) (SEM. NET LUTE	LMT
Transfer						ARIANE MU Others ?	CETA			METEO	T-AVIA	ATLAS (II) LUTE	LMT
Interface Target Structure						ARIANE MU	CETA			METEO TAUM-AV.	ATION	? ¥	? ; ¥
Struct. Generation	-					ARIANE MU	CETA	CETA		METEO T-AVIA			
					×	ARIANE - MU	>	ARIANE (Num. attributes)		METEO T-AVIA			
Morph. Generation		ARIANE CETA								METEO T-AVIA			-
Output Text	~ALL									METEO T-AVIA (1 part.)	-		

# Linguistic "Content" (rather, "INTERPRETATION")

Derivation tree Grammar bound ~ Always projective / Representation tree

 Language bound
 Independently defined
 Better for MT
 (As for Compilers)

### Lexical Information

- Lemmas
- Lexical Units MOVE / MOVER / MOTION / MOVABLE
   1 LU = 1 Derivational Family (MOR / SYN / SEM) Several Meaning Subfamilies (OBSERVE / OBSERVANCE) (OBSERVE / OBSERVATION)

Lexical Functions (à la Mel'chuk)

Not (yet) used in MT

[Power]

Lexical units give good (translational) paraphrasing

+ Annex INFO

- Valencies / Case Frames
- Semantic Features

### Grammatical Information

- Morphosyntactic Info (Gender, Number, Category ...)
- Determination
  - Actualization (Tense / Mode / Aspect .....)
- Syntagmatic Category (1st "Level of Interpretation")

# **Relational Information**

SF	Syntactic Function							
	(More than dependency relation, because relates a word or a group of words to same.)	3						
LR	Logical Relation (ARGO, ARG1, ARG2, ARG12, TRL10, TRL21,)							
SR	Semantic Relation (CAUSE, INSTR, CONSEQ, MEASURE, ACCOMP, LOCAL, TOY, DEST, ORIGIN,)							
	$SR \neq LR$							
	SR difficult to compute on arguments							
	For translation of arguments, LR is enough ( $+$ SEM features).							
	A Key <sup>INSTR</sup> Opens a Door <sup>TOY</sup>							
	One <sup>AGNT</sup> Opens a Door <sup>TOY</sup> With a Key <sup>INST</sup>							
	$A DOOR^2$ Opens (With a Key <sup>INST</sup> )							
	ARGO PRED ARG1							

Levels of Interpretation in Interface Source Structures

Levels MT & al. Systems	K	SF	LR	SR	Sematic Features	Lexical Units	Traces
SYSTRAN	+	+		-	50⊅400↘~100 used	<b>–</b> .	+
LOGOS	, +,	+		?	~100 or more	-	+
CETA	0	0	+	+	~25/30	+	-
METEO	$\bigcirc$	0	·	+	~30/50 (Domain-Spec.)	-	- ?
ARIANE	+	+.	+ 1	+ '	~40/50	+ (Variants Exist)	+
TAUM-AVIA	+	+		+	~100? (Domain-Spec.)	_	+?
METAL	+	+	-	?	~50 (?)	- ?	+
MU	+	+	-	+	~40	-	+

 $\,\odot\,$  Used for analysis but not in interface structure

### Dubious / Incomplete / Ambiguous Analysis

- Almost all systems choose a "Disambiguated" solution, thereby losing the ability to
  - Transfer the Doubt / Ambiguity
  - Treat these phenomena
- Pr. Vauquois' Solution

Use Interface Structures as Objects, & Encode

- Marks of Doubt
- Some Ambiguities / Polysemies
- Indications for later processing

### Examples



For each type of problem, each application defines the default chosen, and the encoding of alternatives.

What about still unmentioned Feature-Structures (complex DAGs)?

2 Extremes

- ① Simply labelled representation trees
- 2 Complex F-Structures on binary derivation trees
  - $1 \rightarrow$  Clumsy handling of set operations or predicates
  - 2 → Clumsy handling of sequences (Order, Repetitions, .....)
- Decorated Trees more "Balanced"
  - Each node has bounded-size Info
  - Natural encoding of set / sequence related Infos
- But Complex F-Structures more explicit

   (Actually, too much for intensive computation: Reentrancy impedes "Divide & Conquer")

### A Simple Example



- 32 -

### A "Rule of Thumb" for using this or that Structure



Of course, some simple Representation Trees / F-Structures may be obtained from Derivation Trees by local attached actions.

But interesting Representation Trees (e.g.  $a^n b^n c^n$ ) usually can not be obtained without a more Complex Device.

A Sketch of the Beginning of a Tentative Approach ..... To Define a Convenient Interface Representation Structure in a MI System

#### Some Motivations

- The main problem of MI is to link speech & NL techniques.
- Symbolic & heuristic approaches have failed in speech processing, while stochastic methods, more "Brute-Force" are (now) doing better.

- The integration through a "Blackboard" seems a good idea, if programmed heuristics (Agenda, .....) are replaced by more parallel, weight-manipulating methods.
- The next problem of MI is basic speed.
  - → Impossibility of using combinatorial approaches with heavy pattern-matching or unification.

### <u>Some Ideas</u>

- Use a working (factorizing) structure all the way through analysis (Blackboard-Like)
- Work by levels, in the way of speech techniques.
- Construct simple structures first, and more elaborate ones on / from the simple ones which get high grades.
- Use vertical activation / lateral inhibition.

# A Possible Macroscopic View



### Content of a Cell

- Main Label & Connexions (?) may be built in the data structure.
   (Not the words, however)
- Simple Attribute Structure (Or direct link to microfeatures)
- Node of Representation Tree
- With a Complex Decoration
- With the correspondences to input (SNODE, STREE)
- & ..... with Complex F-Structure, if needed & computable in real time.

4th Work Session (June 1, 1988)

# Specification and Application of Grammars and Dictionaries

# **Outline of Presentations**

### Grammars & Dictionaries

Static / Dynamic

General / Specific (to MT) but Neutral

/ Specialized (System, Phase)

& their writers

#### Main Computational Methods in NU

Combinatorial/Heuristic Pure/Impure

How to reconcile the 2 main trends?

• Impossibility of deriving automatically an efficient program from a functional or relational specification.

• Techniques from COMP SC : Invariants

• Techniques from AI : Specific Heuristics

A Suggestion for MI Systems.

From Data-Bound to Goal-Bound

&

From Combinatorial to Heuristic

#### **Grammars & Dictionaries**

- Static / Dynamic
   Declarative + Procedural
- General/.....
  - Typically hand-made, dictionaries rarely computerized, grammars never
  - Not usable directly for MT or MI

- But indispensable resources ..... / Specific to MT but neutral
- Ex: MU-Project, CALLIOPE, LOGOS LEXDB + Automatic Generation of "Running" Dict.
- For Grammars: GETA's Static Grammars All other MT Projects have only Dynamic Grammars. ...../Specialized (System, Phase)
- Exact format required eg. MORPH. ANAL.: Access by MORPH, .....
- Use of predefined set of codes (Templates, .....)

### **Necessity to Separate Grammars & Diecionaries**

- Not the same writers, different competences & creation techniques
- Dictionary / Domain
   Grammar / Typology
  - iranninai / rypology

An MT or MI System is <u>always</u> specialized.

..... Success of Expert Systems,

Failure of General Systems .....

Computational Side (Time):

- Dictionary-related phase

(ng)

) (ng log d)

g : grammar d : dictionary

length

n:

- Ň
- Grammar-related phase G = g + d  $\bigcirc (nG) \bigcirc (n^{\begin{cases} 2 \\ 3 \end{cases}}G^2) \bigcirc (n^{6}2^G) \bigcirc (n^3 3^{GG})$ RARE ! EARLEY ..... HGPSG GPSG .....

Contrary to beliefs of theorists, a lexical approach does not guarantee small grammars.

Also

Actual dictionaries proposed by XYZGs proponents are

- Very small
- Very poor (quite poorer than <u>usual</u> MT dictionaries)

Actual grammars proposed by same are also

- Very small
- Too "clean" (no interest for <u>usual</u> phenomena) See Tsujii's, Tomita's Comments.
- It is quite possible to get a grammar larger than 200 p (10000 lines).

### Main Computational Methods

### Combinatorial Methods

 Pure : Produce all solutions CF-Grammars in the sixties AUGM. CF-G in the sixty-fives

(& now - METAL, .....)

U-Grammars in the eighties

("Packing" is the new word for "Factorizing", used-of course-since 1960)

• Impure : Deliver <u>1</u> or <u>some</u> solutions

- Filters old idea

Mel'tchuk & Kulagina (1956) PIAF (1970)

..... All "Sequential" approaches

Preferences

~ 1967-73

CETA System (For Syntax) WILKS (Preference Semantics)

GETA's Systems (Integrated in AS)

- Weights

METAL, Critique

All solutions are computed first, then some choice is effected.

### Heuristic Methods

Try not to compute all possible solutions

- Pure : Search Strategies

   (<u>Without</u> storing of partial results)
   Backtrack Methods (PROLOG, ATN, ROBRA)
   Idea : The first found is the best.
- Impure : Associated to combinatorial base
   ATEF (ND-Finite State + Heuristic "Functions")
   Agenda + Blackboard idea

   (Not yet used in MT, but ..... in MI ?)

### How to reconcile the 2 Main Trends?

- In MT, the 2 types of techniques have achieved success.
- In speech processing (recognition), the Data-Bound, Combinatorial Methods (using weights & not preferences) seem to WIN, the "More Intelligent" to fail.
- Quite the contrapy when it comes to "Deep" understanding, or knowledge-based processes.
- To this day, it is impossible to derive automatically an efficient program from a functional (relational) specification
  - No practical results of program synthesis research
  - No actual use of applicative programming
     ..... Also, incredibly difficult to debug even toy examples
  - Theoretical limitations, eg:

 $I\Psi = \{x \mid \psi_x \mid \Psi\}$  not recursively enumerable

→ Computational linguists or linguistic-oriented computer scientists are not likely to discover this GRAAL

→ Some procedural component (performance related ?) has to be ..... programmed. Just how ?

#### Techniques from Classical Computer Science

- Use the (static) grammars as specification
- Construct (EX NIHILO) some program (in ROBRA, GRADE, ANT, .....)
- Enrich the program with
  - Normal comments explaining the relation to the specification
  - Special comments, of formal nature, the invariants, at well-chosen points
- Enrich the programming language (SLLP) with run-time procedures for checking the truth value of the invariants (always, periodically .....) and take appropriate action

### Techniques from AI

- Declarative Component
- Meta-rule Component

How to solve local problems, or heuristics

- Preferences, weights, interaction with user or knowledge base
- General Strategy

At the highest level, global handling of non-determinism

- Depth / Breadth / Best First
- Development of *n* solutions in parallel
- Freezing / Activating .....

### A Suggestion for MI



5th Work Session (June 22, 1988)

# Specialized Languages for Linguistic Programming

# **Outline of Presentations**

1. 3 Possible Criteria for Classifying SLLPs

- Types of Production Systems
- Rule Mechanisms (Direct, P-M, Unification)
- Control & Modular Organization
- 2. SLLPs of the Direct Type
- 3. Pattern-Matching Based SLLPs
- 4. Unification-Based SLLPs
- 5. Discussion on the place of these 3 types of SLLPs in machine interpretation.

### 3 Possible Criteria for Classifying SLLPs

• Types of Production Systems

Addition - Usually followed by cleaning-up



ex: Q-Systems, Chart Parsers



### Transduction



• Rule Mechanisms

Direct:



No global search in the working structure(s)

② Main element in the left-hand side is <u>constant</u>.

ex: Finite-state transducer (ATEF)



- 44 -

### If variables appear, they are merely instantiaded also: Context-Free Parsers, ATNs

Pattern-Matching : LHS with variables global search But structure has no variables



ex : Q-Systems ROBRA, GRADE

Unification : Variables on both side ex : Logic Grammars (MG, DCG, XG / LMT ....) "Unification-Based" (F-Structure Based) Grammars

Grammars

### • Control & Modular Organization

- Subgrammars
  - Programmed

r:  $A \rightarrow \alpha/G1/G2$ 

Matrix

 $[r^1, r^3, r^{10}]$ 

- Q-Systems
- ATNs, REZO

$$\mathbf{A} \to \mathbf{a}_1 \mid \mathbf{a}_2 \mid \dots \mid \mathbf{a}_n$$



- ATEF, ROBRA, GRADE

a (sub)network corresponds to all rules with same L.H.S.

In reality, no subgrammars

cf. Salomaa

Simple block of rules

Real subgrammar organisation



Q-Systems: Sequence (MTL) Trial & Error (Prague)



Levels of Rules

METAL (with CFG)



For Machine Interpretation,

The Main Criterion is Time Complexity (~ Real-Time Constraint)

Hence Study by the Rule-Mechanism Classification

#### SLLIPs of the Direct Type

(No global search for a pattern, essentially data-driven)

- Regular (L-R) String Transducers, even if non deterministic.
   ex. ATEF (ND) Morphological Analysis
   SYGMOR (D) ------ Generation
- All CF-based formalisms, equipped with reasonable algorithms, if the computation does not entail matching / unification on unbounded structures

# METAL, PNLP (IBM YTH) ATN, REZO

• Tree-Transformational Systems, if in "Anchored" mode ex. TTEDIT (CETA)



### Pattern-Matching Based SLLPs

Global search for a pattern over an entire "Object Structure"

- Tree-Transformational Systems
   ROBRA, GRADE, TELESI (Chanché)
  - → Because they use substitution, they must solve conflicts /

Q-Systems

Not like chart parsers because :

- May create new nodes
- Don't test for equality



R1

R2

R1

These systems are extremely powerful, & usable not only in analysis, but in transfer & generation

Cost: ~ 50 times cost of direct SLLPs (If programmed comparably and if linguistic data comparable)

### Unification-Based SLLP's

All those where unification of potentially unbounded structures is the main operation.

- CF-Based
  - ex. D-PATR, as used in ATR



• & where the result is really a complete attributed tree.

& apparently also not JPSG (Gunji)

- MG (Metamorphosis Grammars based on CSG)
- STCG (or "Static Grammars") Zaharin Here, identification rather than unification (Variable on 1 side may not be substituted)

### Cost of U-Based SLLPs ?

• In the previous sense:

Not yet any <u>actual</u> measurement, on complete systems As opposed to direct & PM-based Preliminary experiments → Staggering cost Theoretical complexity (Book by Barton, Berwick & Ristad) NP-complete or more

But, in practice,

Although some formalisms claim to be theories, thus exposing themselves to this reproach,

All reduce to tools, hence SLLPs

 $\rightarrow$  eg. D-PATR from many U-formalisms

The real problems are :

- The cost of the <u>used</u> unification (Real/Pseudo)
- The added cost if <u>types</u> and <u>other operations</u> are allowed, which is <u>necessary</u> for a SLLP to be useful (numbers, strings, ..... intersection, difference .....)

### Discussion on the place of these 3 types of SLLP in MI

• Recall the considered levels & scope

Real-world, pragmatics	Whole situation, not only one dialogue	$\infty$
Dialogue dynamics (in planning / understanding)	One dialogue	<2000
Dialogue semantics (speech acts, participants)	Dialogue "Window" or total dialogue	<100 ?
Utterance semantics	Utterance + Window	
Utterance (abstract) structure (complex construction)	Utterance (1-n Sentences)	<200
Complex syntagmatic groups	Sentence at most	<80?
Simple syntagmatic groups	No recursive embedding	<15?

Elementary syntagmatic groups

Kernels eg. John Smith Finite State

Words [Syllables] [Pphonemes]

Signal "Frames"

1f = 50ms



- 51 -

How could one "Mix" these different types ?

Idea : Unify the SYNTAX (of the SLLPs) Equip with several "Rule Engines" Unify the Data Structures

Data Structures

There should be

Nodes with Label, {Attributes}, [F-Structure]

Trees (ordered)

(Loopfree) Graphs Charts? Q-Graphs? Lattices?



1) First construct only

[----] Skeleton

eg. by attached transformations

2) Then make the <u>selected</u> representation(s) more explicit by computing the F-STRUCT 3) Suppose weights have been used in the construction.

We get a natural way to translate the "Important" parts (of TSUJII) at a higher level than the rest.



Result of 1+2 3: PRVNE

Note : "Important" does not mean "domain-related" any more, just dependent on ..... How good <u>we</u> judge our analysis.

# Translated at "Linguistic" Level



Translated at "Knowledge" Level (From the Semantic Interpretations)

Note: Translation at linguistic level does not have to be as bad (word for word) as in TSUJII's sketch !

Lattices rather than Q-Graph or Charts?



Lattices seem more used in speech processing.

# Organization of the Whole MT / MI System

# **Outline of Presentations**

### I. The Importance of Size

- Quantum Leaps in Methods
- Size & Time Complexities
- Mockups, Prototypes, Operational Systems

### $\Pi$ . Environments

¢	Preparation	:	Lingware Versions, Tests Test Corpuses
0	Experimentation	•	Compiled Packages, Measures Validation Corpuses
•	Debugging	*	Variety of Levels, Tools
0	Production	•	End-User Interface Translation Request, Revision .

III. Some Useful Concepts for Modularity

- Lingware : Steps, Phases, Components, Versions, .....
- Texts : Corpuses, Transcriptions, Translations, .....

#### The Importance of Size

• Quantum Leaps in Methods

What works on small, closed problems is usually not extendable.

- e.g. Morphological analysis may be unnecessary if you work on 20 sentences.
  - Trigram methods may work for 1000 words (10<sup>9</sup> ELTS), unlikely for 100,000 (10<sup>15</sup> ELTS).

Usual steps in MT development

		Man  imes Year	Dict. Size
20 sentences	(1 page, 250 w.)	$1-3 \times 1$	200 w.
200 sentences	(10 pages, 2500 w.)	$3-4 \times 3$	2000 w.
2000 sentences	(100 pages, 25000 w.)	$\sim 10 \times 3$	20000 w.
20000 sentences	(1000 pages, 250000 w.)	$\sim 20 \times 3$	~100000 w.

Problems in grammar & dictionary management may force to abandon methods successful until a certain point.

- e.g. CF-based grammars begin to be unmanageable with a few hundred rules.
  - Direct dictionary coding OK for <10-15000 w, by then a lexical data base management becomes necessary.
  - Although example of complex transfer dictionary
- Size & Time Complexities

SLLPSs should be concise enough

Remember Charniak's description of "Paint" (half done !) on 9 pages of "Planner"

The same goes for the descriptions of units of translation

- For this reason, trees are preferred over graphs (note that any graph can be covered by trees). In any case locality is important.
- Linguistic trees should not have too many or too few nodes.



Compilation & execution time crucial

The best formalism is useless if it takes an hour to compile a small grammar and a day to translate 10 words (eg. EUROTRA).

- Mockups, Prototypes, Operational Systems
  - Mockup: Only the heart (1, 2, 3 SLLPSs) No environment Usually 1 language-pair & 1 version
  - Limit : Something like METEO during 1st year of development at TAUM.

Prototypes

Small

Debugging environment DB management of lingware Almost nothing for the texts Nothing for communication / operationa

e.g. -

D-PATR Only 1 text per grammar

Non sharable

(But D-PATR is far from a prototype environment for MT - no transfer, generation)

- LUTE ?

Large

dictionaries

And lots of functions on texts

But no integration in Documentation System

DBMS adequate for several versions, large

Network

e.g. ARIAVE (-78, -85) cf. COLING-82, JACL 86/1 MU

**Operational Systems** 

End-User interfaces

Bilingual editor Command panels (Easy Dict. System)

e.g. SYSTRAN-Cee (+ TERMEX/MINITEL .....) LOGOS

# ATLAS (Fujitsu), PIVOT (NEC), HICAT (Hitachi), AS-TRANSAC (Toshiba)

### Environments

- The environment ..... of the environment
  - Dialogue language
  - Dialogue style (abbreviated, long, detailed)
  - Source language code
  - Target language code
  - Corpus name
  - > Trace parameters for each phase
- Preparation

...

- > Editor settings for various types of
  - [Linguistic data (free / fixed format, ....)

[Texts (input transcription)]

- > Automatic compilation (partial / total)
- Lists (with all possible sorting orders)
- Cross-references
- Synthetic views across files
   (e.g. Produce all information for a set of lexical units)
  - small tests
- Experimentation
  - Facility for producing compiled packages (memory images), quickly loaded / executed
  - > Possibility to trace on one or several validation corpuses,  $\pm$  selection
  - Storing of intermediate results to
    - Save future computation
    - Modify them (tree editor, frame editor) to produce normalized test inputs for subsequent phases.

Debugging

The facilities of the implementation language are not adequate

- $\rightarrow$  Various views of the unit of translation input/output of a phase
- $\rightarrow$  Various trace / step-by-step parameters
- → If the units of translation are large, it is also important to easily isolate the erroneous part, in order to deduce the size of the structures produced.
- Production

Important to be able to start the translation of several brochures (say, 1 corpus per brochure & 1 text per section / abstract) in less than 1 minute, when rushing out at the end of the day .....

- → In a LAB, the production environment should offer the basic functions with which the operational system(s) will be built.
- e.g. ARIANE as background for CALLIOPE



Luxurious Text Processor

# Modularity, ma jolie ..... One (ARIANE) Space

- Texts
- Linguistic Applications
- Linguistic Modules (memory images)
- (Lisp) Programs  $\rightarrow$  You may prefer C or PROLOG !
- Natural On-Line Dictionaries

- System Files (e.g. Internal ↔ External Names)
   (e.g. Values of environment variables)
- Lexical Data-Base, from which MT Dictionaries are constructed (e.g. MORPH. ANALYSIS, GENER)
- Transcriptors (in LT)
   e.g. 'Déjà là, DUPONT ?' → '\*DE!1JA!2 LA!2, \*\* DUPONT'

• Personal files of individual user.

Lingware

Step

Analysis Transfer Generation one could add Speech recognition Speech synthesis

Understanding

Phase

Autonomous part of a step, written in <u>1</u> SLLPe.g.Morphological AnalysisAMStructural AnalysisAS

Component

Element of a phase

e.g. Grammar G1, G2, ..... Dictionary D1, D2, .....

Selection of components gives possibly several realisations (G1 + D2)



# Precedence Graphs for Compilation of Components



### Textware

Texts of same language, typology, domain, transcription Corpus (imagine ASCII/JIS)

Not only 1 file, but complex file set : Text

Source Text

Brute	As inputted, in which format			
Revised	Often, we have to preprocess slightly			
•	(spelling mistakes, OCR errors)			
May be in more than 1 format,				
	e g ISOØ25 (French Keyboard) & SCRIPT			

с.д. JIS & JTEX

CARIX & TETHYS (implementation 30%)

Descriptor

Gives logical structure (sections, paragraphs, sentences) used to segment in units of translation

Figures, Formulas, .....

Formats (HORS-TEXTE) Brute / Revised

Translation	(Brute/Revised)
	1 for each $A + T + G$ realization
Intermediate	1 for each variant of each phase
results	•

### Tools on Texts

List of Words	3	On 1 on day	On 1 on governal		Texts		
Concordance	s	OII I OF SE	On 1 or several				
Unknown Words (± Context)							
Merging / Splitting							
Printing	± Tran ± Rev	nslations isions	}	Problem	n of onization		

Editing

usually, by segments

Morphosyntactic Search

- e.g. Look for Art Adj Adj Noun
- $\rightarrow$  Stochastic techniques have been used (Spirit System) successfully See also Lancaster group.
- More powerful techniques for "Structural Search" should be reveloped.

See DEREDEC (Montreal) as first example.

• The Example of ARIANE-85



• Component "Matrices"



TRACOMPL (for Transformation & Complement of DECORATIONS)

always 1 fik only

DV

- 65 -



& dont't forget to Backup !!!

& Levels

In user space (DUPLG / DESTRUL COMMANDS) Use Backup ("Fixed") user space (TRGLGVM) And tapes ..... (TRGLGVM COMMAND)

### What else?

Many things not in ARIANE, ideas could come from various sources :

• Integrate everything under the editor / window system, as in lute, instead of calling the editor on speicific files.

Integrate indexing workbench, with

- ± Lexical Data Base
- + Indexing Aids
- Access to "Natural" Dictionaries (MINITEL .....)

Reconsider the organization, to distribute 1 user space over a (physical or logical) network, with several users in parallel.

FRAMES

For machine interpretation

Add several forms of texts

PHONEME	Lattice
SYLLABLE	or
WORD (phon	Grid
or orth.)	

Tools for stochastic classification

Manual tagging of samples

Construction of HHM ("Learning") etc.

Finally, the structure of a whole MT System is like that of a software factory.

• A useful partition of sub-environments is by users :

	End-User	should come FIRST !
	Linguist	)
	Lexicographer	
+	Phonologist	Developers
+	Cogniticist(!)	J. And S.
	Computer Specialist	Special functions, tools,

• Perhaps the <u>execution</u> can also be distributed on several machines (cf. HEARSAY), the environment should then take it into account from the start.

'n.