

# WP2 - Semantic Parsing and Generation of Documents and Documents Components

*Claire GARDENT, Bikash GYAWALI,  
Anastasia SHIMORINA*

**CNRS / LORIA**

*Samuel CRUZ-LARA*

**University of Lorraine / LORIA**

## WP2



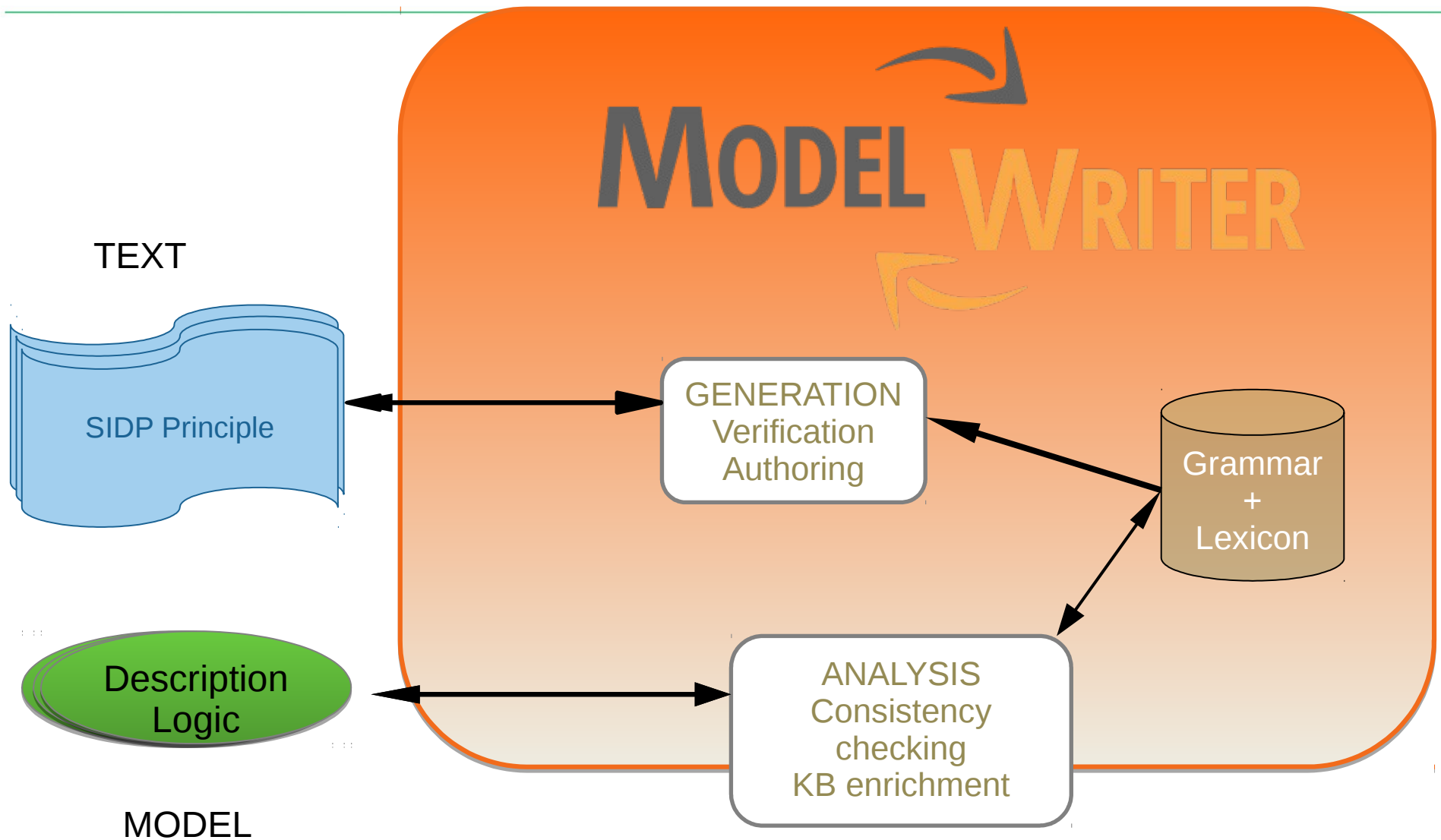
Goal: Provide tools and methods for:

- Converting texts to models and models to text
- Annotating text fragments with model elements

Tasks:

- T2.1 Data Collection
- T2.2 Semantic Parsing
- T2.3 Natural Language Generation
- T2.4 Definition of a common target semantic language
- T2.5 Development of a Semantic Parser and of a Natural Language Generator

# Synchronising Text and Model



## Related Work on Ontology Learning

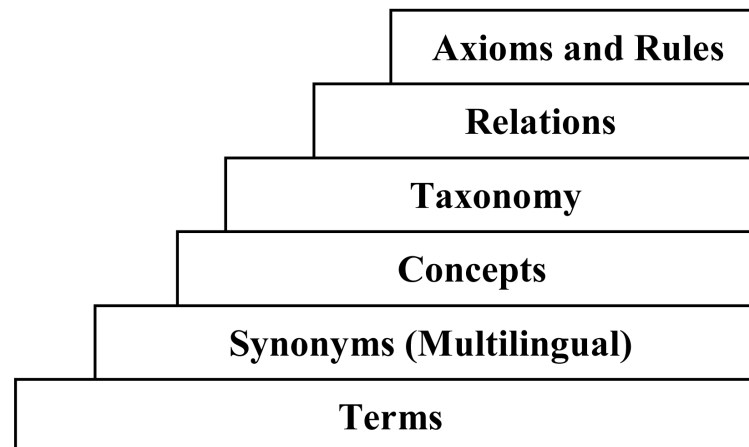
---

### Ontology Learning from Text

[Mädche and Staab 2000, Volker et al. 2007, Tablan et al. 2006, Zouaq and Nkambou 2008]

*Restricted expressivity. Not applied to sentences (complex axioms).*

*Limited to : definition of new classes, creation of hierarchies between classes, definition of object and data-type properties, creation of instances, and setting of property values for instances*



## Related Work on Semantic Parsing

---

### Deep parsing

[Currant et al. 2007, McCartney and Manning 2007 ]

*First Order Logic Representations close to initial text. Trained on newspaper text (Penn Tree Bank).*

*==> Not easily adaptable to Description Logic and SIDP text.*

### Domain Specific Semantic Parsing

[Ge and Mooney 2009, Wong and Mooney 2007]

*==> Require parallel text-data training corpus.*

### Open Domain Semantic parsing

[Kwiatkowski et al 2010, Bordes et al. 2012, Kwiatkowski et al 2013, Berant et al., 2013, Bordes et al. 2014, Wang et al. 2015]

*==> Restricted to questions. Require parallel question-answer training corpus.*

## Related Work on Text Generation

---

### **Symbolic Approaches**

[Dimitrios et al. 2007, Androtsopoulos et al. 2013, Power et al. 2010, Bontcheva et al. 2004 ]

*Heavily dependent on hand-written modules.*

### **Machine Learning Approaches**

[Wong et al. 2007, Belz 2008, Angeli et al. 2010, Chen et al. 2008, Konstas and Lapata 2012a and 2012b]

*Require parallel data-text training corpus.*

### **Pattern-Based**

[Duma et al. 2010, Blake et al. 2013, Schilder et al. 2013]

*Require large quantity of parallel or comparable text-data training corpus.  
Limited Semantic Variability (set of RDF triples).*

## Reversible Processing: Text <--> Model

Analysis: SIDP Rule  $\rightarrow$  Description Logic

Generation: Description Logic  $\rightarrow$  SIDP Rule

Verification by generation

## Semantic Parsing of Complex Axioms

*Pipe shall be identified by labels*

*$\text{Pipe} \sqsubseteq \exists \text{identificationArg2}^-. (\text{Identification} \sqcap \exists \text{identificationArg3} . \text{Label})$*

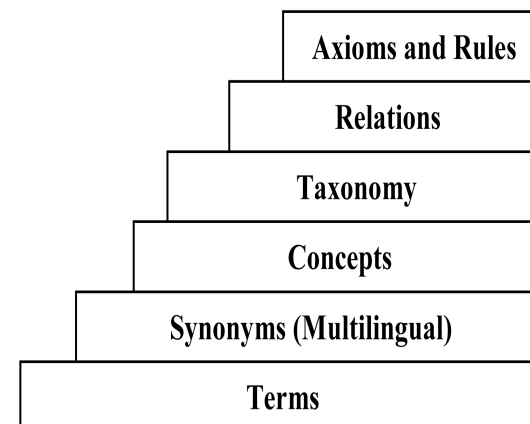
## Executable Semantic Parsing on DL KBs

The output of semantic parsing is used to update a Description logic Knowledge Base and check the consistency of SIDPs (system installation design principle)

## Genericity

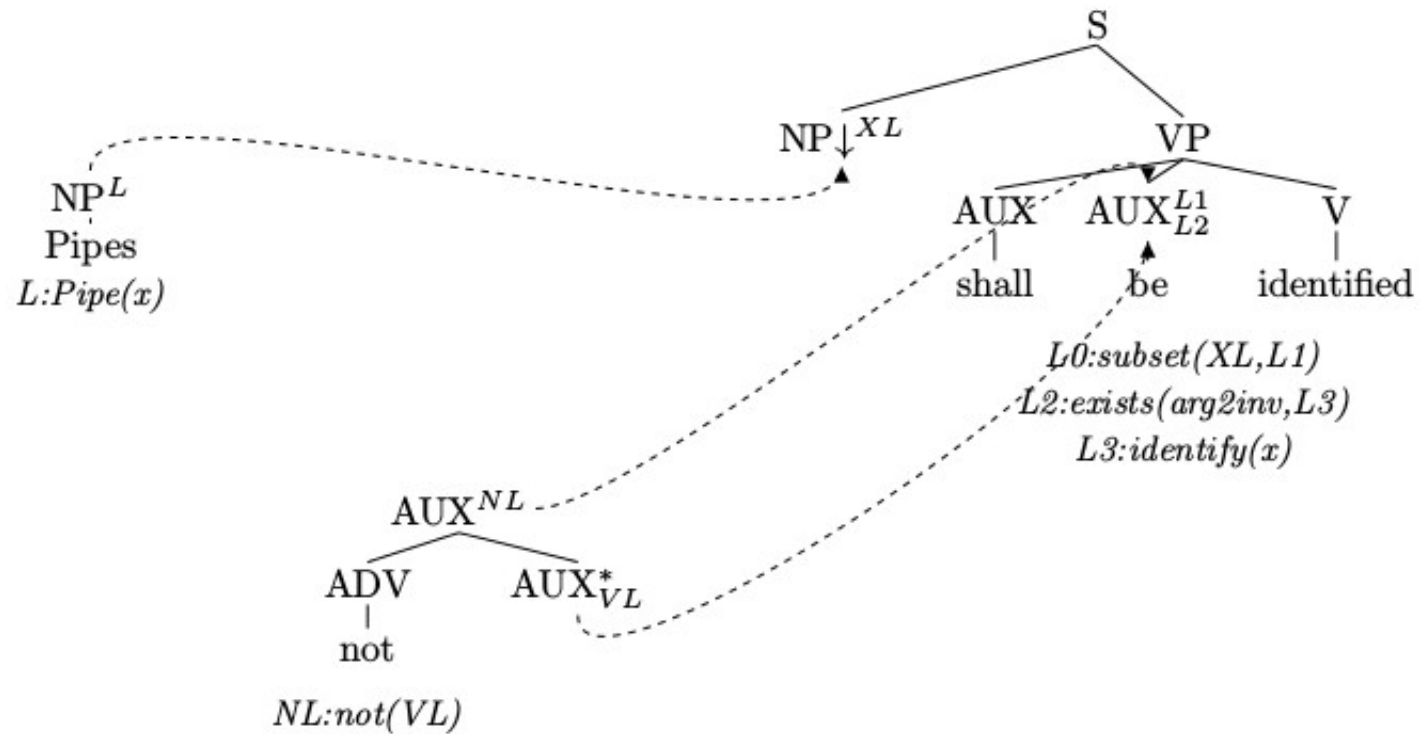
No training corpus required.

Adaptation to a new domain through grammar adaptation, extension or induction



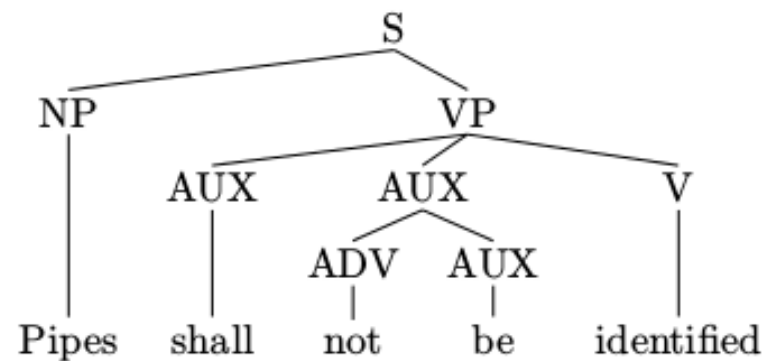


# Grammar-Based Parsing and Generation





# Grammar-Based Parsing and Generation



$PL: Pipe(x)$   $L0: subset(PL, NL)$   $NL: not(VPL)$   
 $VPL: exists(identifyA2inv, VL)$   $VL: Identify(x)$

$Pipe \sqsubseteq \exists \neg identifyA2^{-}.(Identify)$

# Semantic Variations

Logical Operators	
Only S shall be used by O	$\neg S \sqsubseteq \neg \exists use A2^-. (use \sqcap \exists by.O)$
S should be used by <u>all</u> O	$O \sqsubseteq \exists by^-. (Use \sqcap \exists use A2.S)$
S shall <u>not</u> be used by O	$S \sqsubseteq \neg \exists use A2^-. (Use \sqcap \exists by.O)$
Word Order	
S shall be used by O <u>only</u>	$S \sqsubseteq \neg \exists use A2^-. (Use \sqcap \exists by. \neg O)$
<u>All</u> S shall be used by O	$S \sqsubseteq \exists use A2^-. (Use \sqcap \exists by.O)$
Arity	
S shall be used	$S \sqsubseteq \exists use A2^-. (use)$
S shall be used by O	$S \sqsubseteq \exists use A2^-. (use \sqcap \exists by.O)$
S shall be used by O on PO	$S \sqsubseteq \exists use A2^-. (use \sqcap \exists by.O \sqcap \exists on.PO)$
Sentence Structure	
S shall be used by O <u>before</u> entering connections	$(Use \sqcap \exists use A2.S \sqcap \exists by.O) \sqsubseteq \exists before.(Enter \sqcap \exists enter A2.Connections)$
Modifiers	
S shall be used <u>directly</u> by O	$S \sqsubseteq \exists use A2^-. (Use \sqcap \exists directly.(\exists by.O))$
S shall be used by O <u>between</u> C and D	$S \sqsubseteq \exists use A2^-. (Use \sqcap \exists use A3.(O \sqcap \exists between A1^-. (Between \sqcap \exists between A2.C \sqcap \exists between A3.D))))$

The Semantic Representations output by the parser are converted to Description Logic

$l_0 : A(x)$	$\Rightarrow :A$
$l_0 : \text{exists}(R, l_1) \quad l_1 : C$	$\Rightarrow \text{ObjectSomeValuesFrom}(:R \ \tau(C))$
$l_0 : \text{subset}(l_1, l_2) \quad l_1 : C_1 \quad l_2 : C_2$	$\Rightarrow \text{SubClassOf}(\tau(C_1) \ \tau(C_2))$
$l_0 : \text{and}(l_1, l_2) \quad l_1 : C_1 \quad l_2 : C_2$	$\Rightarrow \text{ObjectIntersectionOf}(\tau(C_1) \ \tau(C_2))$
$l_0 : \text{not}(l_1) \quad l_1 : C$	$\Rightarrow \text{not}(\tau(C))$

E.g.,

Pipes should be identified by labels

$l_1 : \text{Pipe}(x) \quad l_0 : \text{subset}(l_1, l_2) \quad l_2 : \text{exists}(\text{identifyA2inv}, l_3) \quad l_3 : \text{and}(l_4, l_5)$   
 $l_4 : \text{Identify}(z) \quad l_5 : \text{exists}(\text{by}, l_6) \quad l_6 : \text{Label}(y)$

`SubClassOf(Pipe ObjectSomeValuesFrom(identifyA2inv  
ObjectIntersectionOf(Identify ObjectSomeValuesFrom(by Label))))`

The formula is added to the AIRBUS KB and Hermit is used to check for consistency

# Experimental Setup and Results for Parsing

---

Grammar: 52 trees

Lexicon: 10781 lexical entries

Parsing algorithm: CKY + Robustness mechanism to skip unknown words

Input: 991 System Installation Design Principles

	Complete Parse	Partial Parse	Failure
Simple SIDP	132	329	24
Complex SIDP	0	496	10
All SIDP	132 (13%)	825 (83%)	34 (3%)

# Updating the Model using Parsing Results (Complete Parses)

CONCEPTS	Nb. of new Concepts	184
	Nb. of Existing Concepts	30
PROPERTIES	Nb. of New Properties	62
	Nb. of SIDP Axioms (from Parsing)	132
SIDP AXIOMS	Nb of Invalid Axioms	0
	Nb. of Redundant Axioms	2
ALL	Total Nb. Of Added Elements	376
	Nb. of Axioms in Initial KB	12469
	Nb. of Axioms in Enriched KB	13029

# Updating the Model using Parsing Results (All Parses)

CONCEPTS	Nb. of new Concepts	667
	Nb. of Existing Concepts	79
PROPERTIES	Nb. of New Properties	98
	Nb. of SIDP Axioms (from Parsing)	957
SIDP AXIOMS	Nb of Invalid Axioms	61
	Nb. of Redundant Axioms	125
ALL	Nb. of Inconsistent Axioms	20
	Nb. of added SIDP Axioms	749
	Total Nb. Of Added Elements	1514
	Nb. of Axioms in Initial KB	12469
	Nb. of Axioms in Enriched KB	14650

## Generation Results

---

Grammar: 52 trees

Lexicon: 10781 lexical entries

Generation algorithm: Tabular + Polarity Filtering

Input: 957 Description Logic Axioms derived from the AIRBUS System  
Installation Design Principles

	Success	Failure
Simple SIDP	448	13
Complex SIDP	470	26
All SIDP	918 (96%)	39 (4%)



## Verifying Parsing Results Using Generation

BLEU	< 0.33	> 0.32 and < 0.67	> 0.66
Complete Parses (S)	1	0	131 (14%)
Complete Parses (C)			
Partial (Simple)	143	117	69
Partial (Complex)	396	91	9
All Parses	540 (56%)	208 (22%)	209 (22%)

Regenerating from the DL formula derived through parsing from an SIDP:

- Produces a sentence identical to the input SIDP for complete parses
- Produces a sentence highly similar to the input SIDP in 44% of the cases for partial parses

## Perspectives and Future Work

---

- Improve lexicon construction using chunking
- Improve coverage on complex sentences including conditions
- Querying the KB (support for AIRBUS engineers)
- Improve robustness and genericity (experiment with deep learning approaches using data expansion techniques and sequence to sequence models)