

Représentation orientée-objet d'expressions idiomatiques dans une méta-grammaire

Sarah Pollet & Agata Savary & Emmanuel Schang & Anais Lefeuvre-Halftermeyer

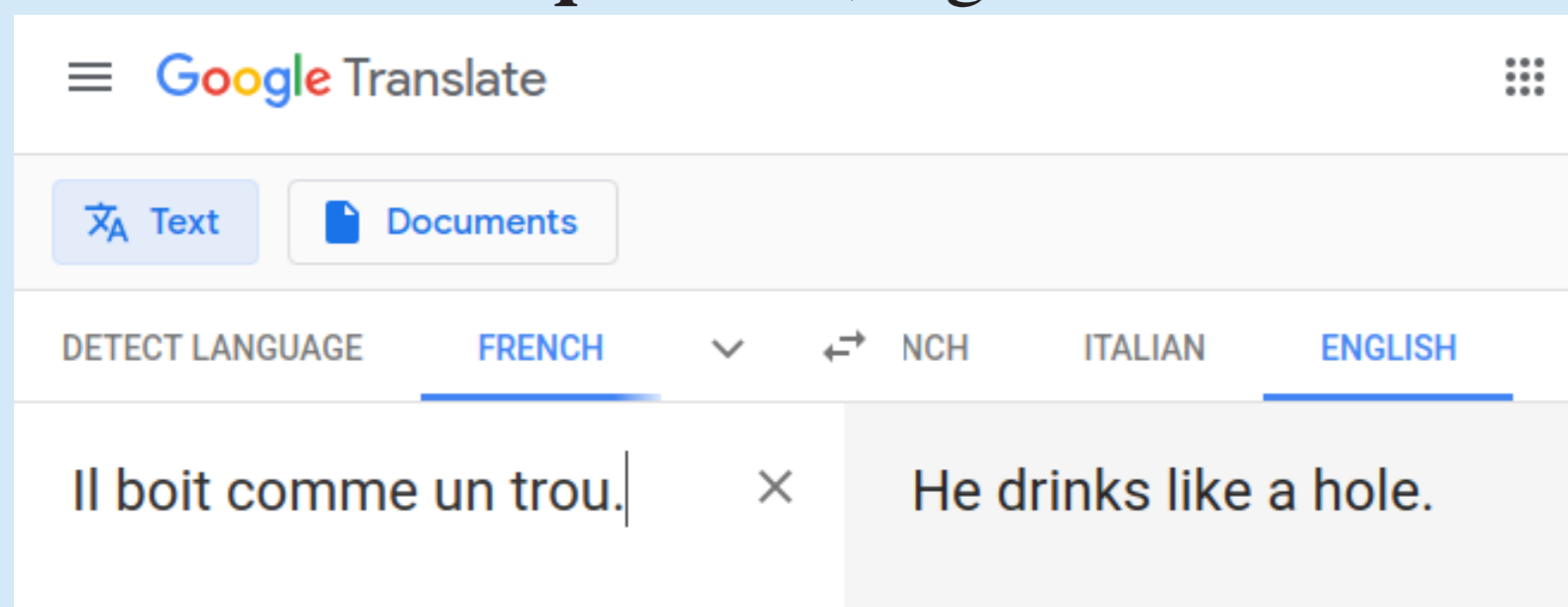
LIFAT, Blois & LIFO, Orléans & LLL, Orléans



Expressions polylexicales (EP: multiword expressions, MWEs)

Si vous **avez** tant **besoin** de **couper l'herbe sous le pied** de quelqu'un, je vous proposerais de **vous en prendre** au **rédacteur-en-chef**, Monsieur **Jean-Marc Petit**.

- **Définition:** Combinaisons de **plusieurs mots** qui possèdent des **propriétés irrégulières** au niveau du lexique, de la grammaire, de la sémantique, etc.
- **Sémantique non-compositionnelle:** Le sens global n'est pas déductible de manière régulière à partir des sens des composants, et des liens syntaxiques qui les relient.
 - ▷ *couper l'herbe sous le pied de quelqu'un* 'empêcher quelqu'un de réussir'
 - ▷ *s'en prendre à quelqu'un* 'prendre quelqu'un pour cible, lui attribuer une faute'
- **Fréquence:** Environ **20%** des mots d'un corpus français de référence appartiennent aux EP (exemple plus haut: 57%).
- **Non-compositionnalité:**
 - ▷ Méthodes informatiques sont **compositionnelles**
 - Phénomènes complexes sont décomposés en problèmes plus simples.
 - Sous problèmes reçoivent des solutions autonomes, qui sont ensuite composées pour fournir solutions globales.
 - ▷ EP sont **sémantiquement non-compositionnelles**, donc posent problème pour les tâches du **TAL** orientées **sémantiquement**, e.g. traduction automatique.



Figement

Une EP est **moins flexible** (variable) qu'une construction régulière de la même structure syntaxique.

Structure	Construction régulière	Expression polylexicale
N Adj	<i>livre bleu</i>	<i>cordon bleu</i> 'excellent cuisinier'
V Det N	<i>manger des salades</i> <i>faire la soupe</i>	<i>raconter des salades</i> 'mentir' <i>faire la tête</i> 'bouder'
V Prep Det N	<i>être sur son siège</i> <i>cuisiner au vinaigre</i>	<i>être sur son trente-et-un</i> 'être très élégant' <i>tourner au vinaigre</i> 's'orienter vers la dispute'
Adj Coord Adj	<i>rouge et bleu</i>	<i>noir et blanc</i> 'nuances de gris'
Det N Aux V	<i>les patates sont cuites</i>	<i>les carottes sont cuites</i> 'la situation est compromise'

Construction régulière	Expression polylexicale	Propriété
<i>livre bleu</i> ≈ <i>livre cyan</i> ≈ <i>bouquin bleu</i>	<i>cordon bleu</i> vs. <i>#cordon cyan</i> vs. <i>#corde bleue</i>	Figement lexical
<i>manger des salades</i> ≈ <i>manger une salade</i>	<i>raconter des salades</i> vs. <i>#raconter une salade</i>	Figement morphologique
<i>elle est sur son siège</i> ≈ <i>elle est sur mon siège</i>	<i>elle est sur son trente-et-un</i> vs. <i>#elle est sur mon trente-et-un</i>	Figement morphosyntaxique
<i>cuisiner au vinaigre</i> ≈ <i>il a fait la soupe</i> ≈ <i>la soupe a été faite par lui</i>	<i>tourner au vinaigre</i> vs. <i>#tourner au vinaigre de vin</i> <i>il a fait la tête</i> vs. <i>#la tête a été faite par lui</i>	Figement syntaxique
<i>les patates sont cuites</i> ≈ <i>nous avons cuit les patates</i>	<i>les carottes sont cuites</i> vs. <i>#on a cuit les carottes</i>	
<i>un texte en rouge et bleu</i> ≈ <i>un texte en bleu et rouge</i>	<i>une photo en noir et blanc</i> vs. <i>#une photo en blanc et noir</i>	

≈: glissement de sens prévisible à partir du remplacement lexical

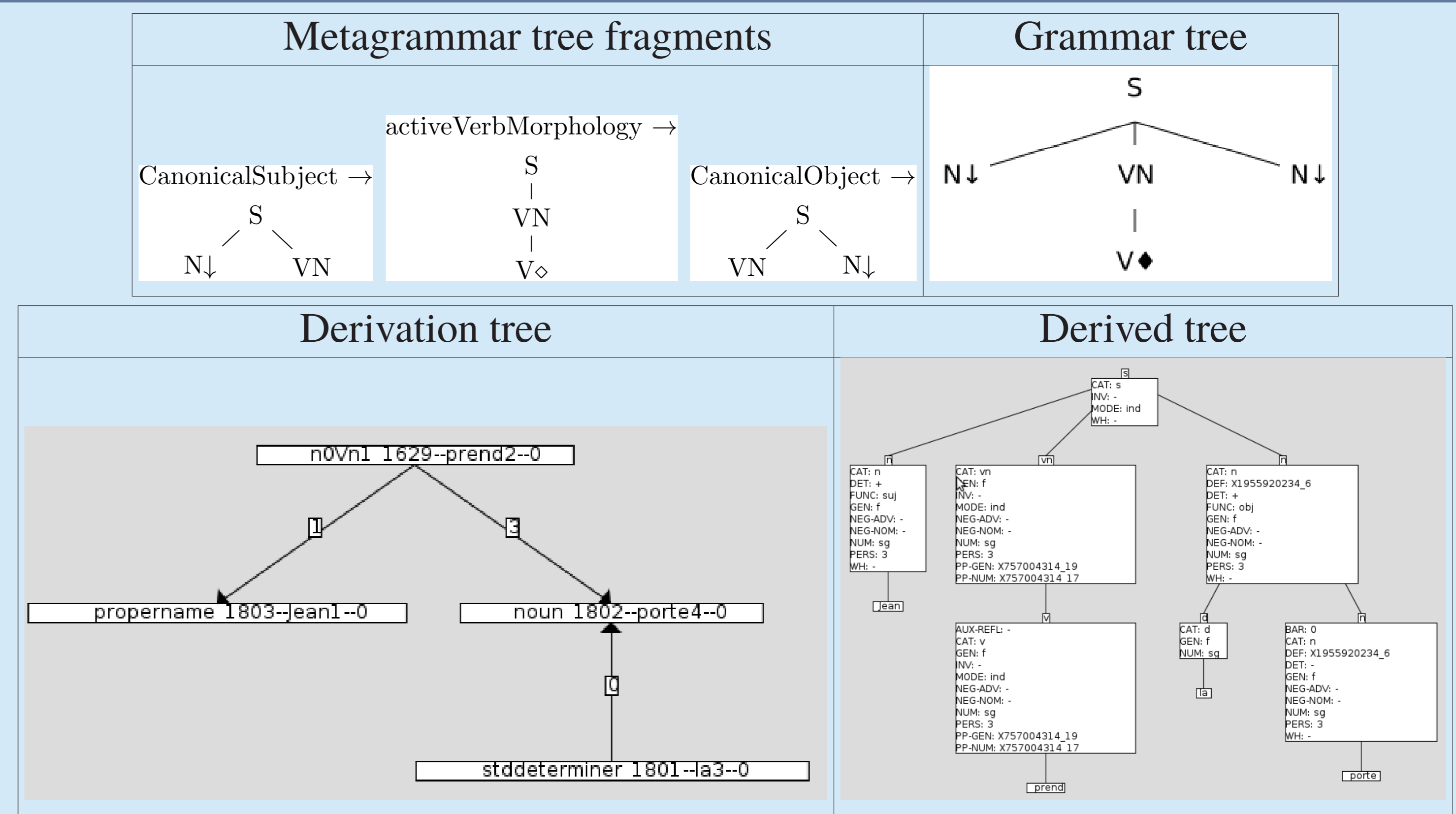
Le figement des EP est une question d'**échelle** et non pas une propriété binaire [Gross(1986), Gross(1988)]

$N_0V (DetN)_1$ expr.	Free subject	Free verb	Free object	Verb red.	Verb infl.	Noun infl.	Noun modif.	Passive	Det. altern.
<i>N prend la pomme</i>	✓	✓	✓	✓	✓	✓	✓	✓	✓
<i>N prend une décision</i>	✓			✓	✓	✓	✓	✓	✓
<i>N tourne la page</i>	✓				✓	?	✓	✓	?
<i>N prend la porte</i>	✓				✓				

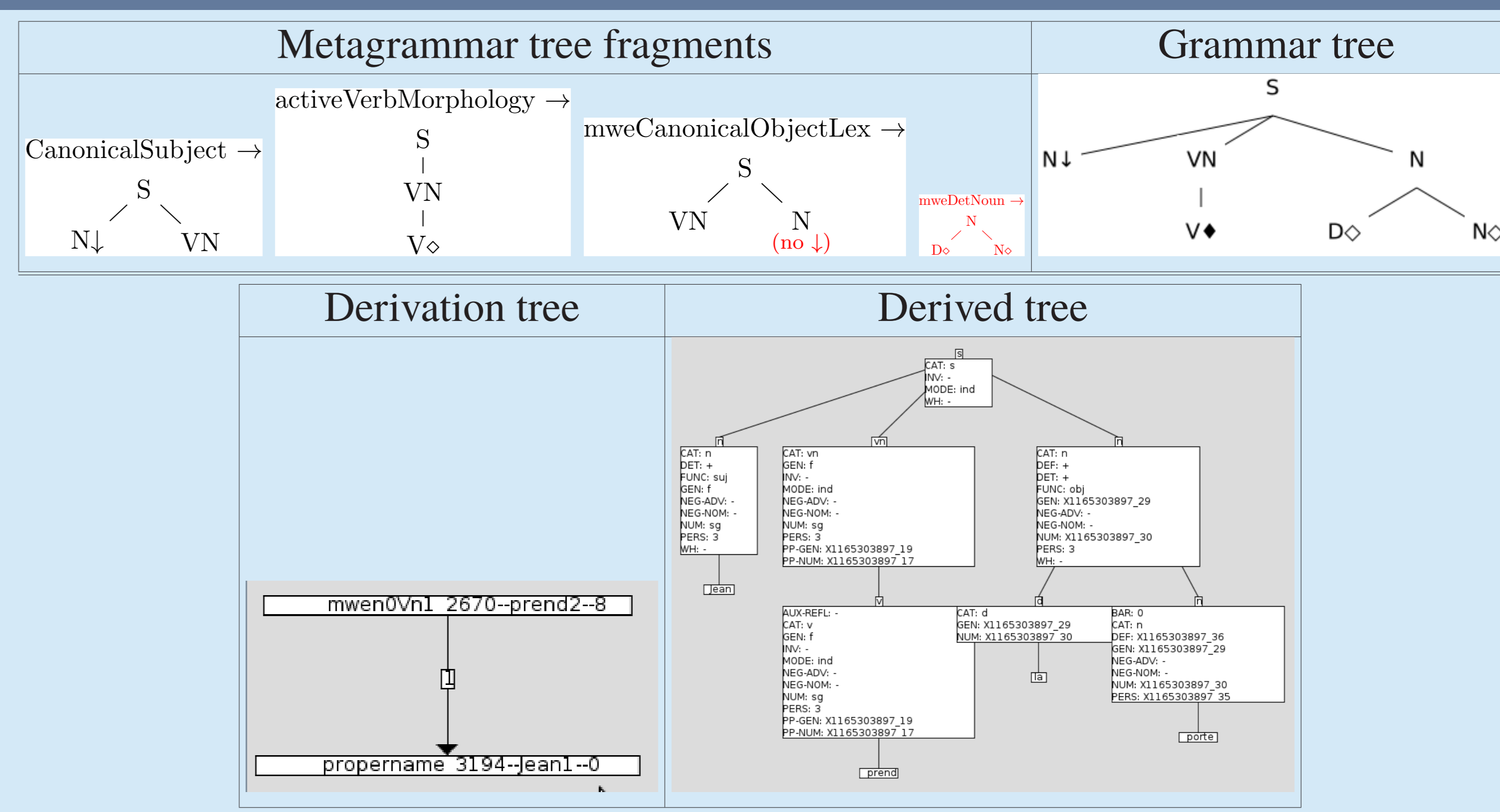
Metagrammatical description

- **XMG** [Crabbé et al.(2013), Petitjean et al.(2016)]
 - ▷ a language
 - **object-oriented** – objects, classes, inheritance
 - declarative – grammaticality = **constraints** rather than procedures
 - notationally **expressive** - modularity, inheritance, conjunction/disjunction of tree fragments, namespaces
 - **extensible** to new dimensions (semantics, frames etc.), formalisms, ...
 - ▷ a metagrammar **compiler** (for each target language, here FS-LTAG) – constraint solver: produces minimal tree models respecting the constraints
 - **FrenchTAG** – French XMG metagrammar [Crabbé et al.(2013)], XMG implementation of the syntactic Tree Adjoining Grammar (TAG) of French [Abeillé(2002)]
 - ▷ 285 XMG classes, 96 families, compiled into 9045 TAG trees
 - ▷ toy lexicon of 555 lexemes, including 248 verbs
 - ▷ covers **literal** readings (by compositionality) but **not idiomatic** readings

Parsing a literal reading: (Jean prend la porte)



Parsing an idiomatic reading: (Jean prend la porte)



Lexicon: morphology and lemmas

```

class Jean
{
  <morpho> {
    morph <- "Jean";
    lemma <- "jean";
    cat <- n;
  }
}

class prend
{
  <morpho> {
    morph <- "prend";
    lemma <- "prendre";
    cat <- v;
  }
}

class la
{
  <morpho> {
    morph <- "la";
    lemma <- "le";
    cat <- d;
    gen <- f;
  }
}

class porte
{
  <morpho> {
    morph <- "porte";
    lemma <- "porte";
    cat <- n;
  }
}

class LemmeJean
{
  <lemma> {
    entry <- "jean";
    cat <- n;
    fam <- propername;
  }
}

class LemmePrendre
{
  <lemma> {
    entry <- "prendre";
    cat <- v;
    fam <- n0vn1;
  }
}

class LemmeLe
{
  <lemma> {
    entry <- "le";
    cat <- d;
    fam <- stddeterminer;
  }
}

class LemmePorte
{
  <lemma> {
    entry <- "porte";
    cat <- n;
    fam <- noun;
  }
}

class mweLemmePrendreLaPorte
{
  <lemma> {
    entry <- "prendre";
    cat <- v;
    fam <- mwen0vn1;
    filter dia = active;
    filter subj = free;
    filter obj = lexicalized;
    filter objstruct = lexDetLexN;
    coanchor ObjDetNode -> "la"/d;
    coanchor ObjNode -> "porte"/n;
    equation ObjNode -> gen=f;
    equation ObjNode -> num=sg;
    equation ObjNode -> modifiable=-;
    filter objtype = canonical;
  }
}
    
```

ICVL internship: Developing a real size lexicon for FrenchTAG

- Original FrenchTAG lexicon: toy lexicon of 555 lexemes, including 248 verbs
- French lexical resources:
 - ▷ **Lefff** - Lexique des Formes Fléchies du Français [Sagot, 2010]: 110 000 lemmas, 520 000 formes
- **Lexicon-Grammar** [Gross 1975, Tolone 2011] - 40,000 EP
- Objectives: convert Lefff and the Lexicon-Grammar to XMG
- Outcome:
 - ▷ XMG lexicon with **112,000 lemmas** and their inflected forms