

# OPEN INFORMATION EXTRACTOR FOR SEVA

**Jitin Krishnan**

Mentor: Patrick Coronado

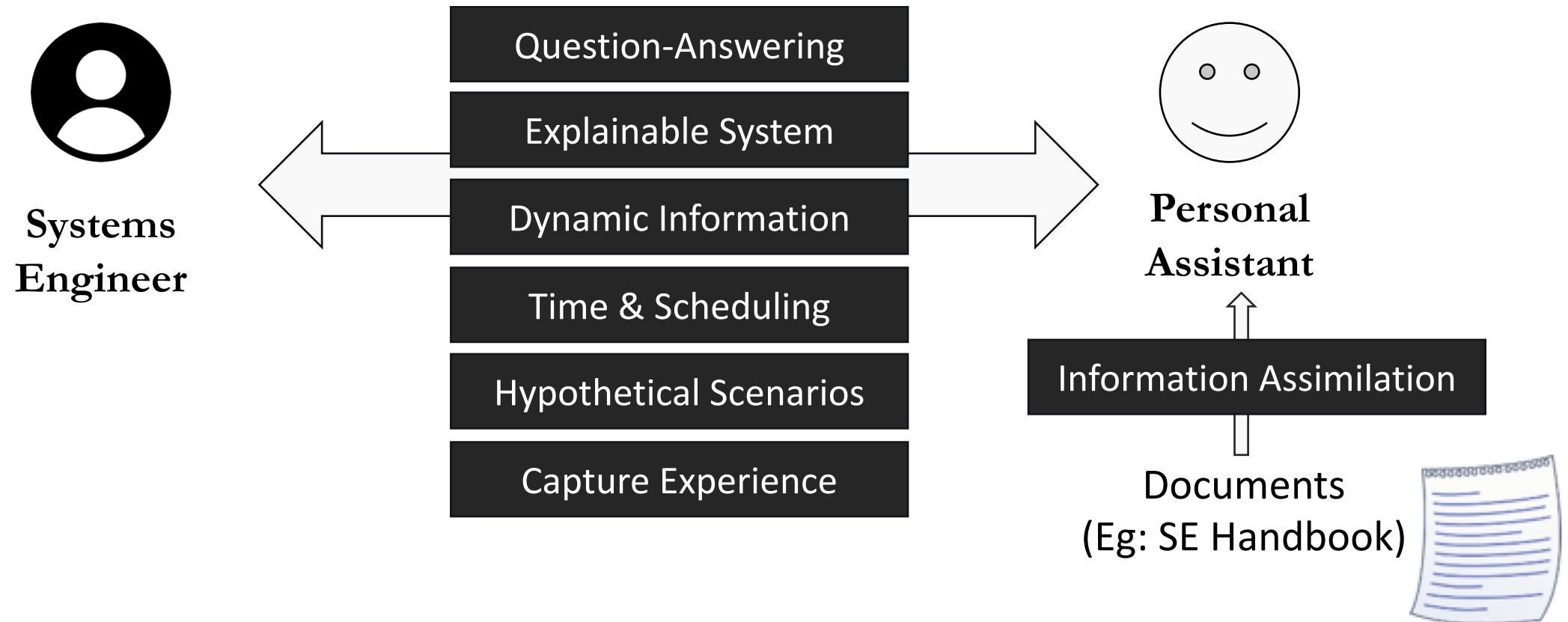
**Other Project Members:**

Trevor Reed (NASA JPL)

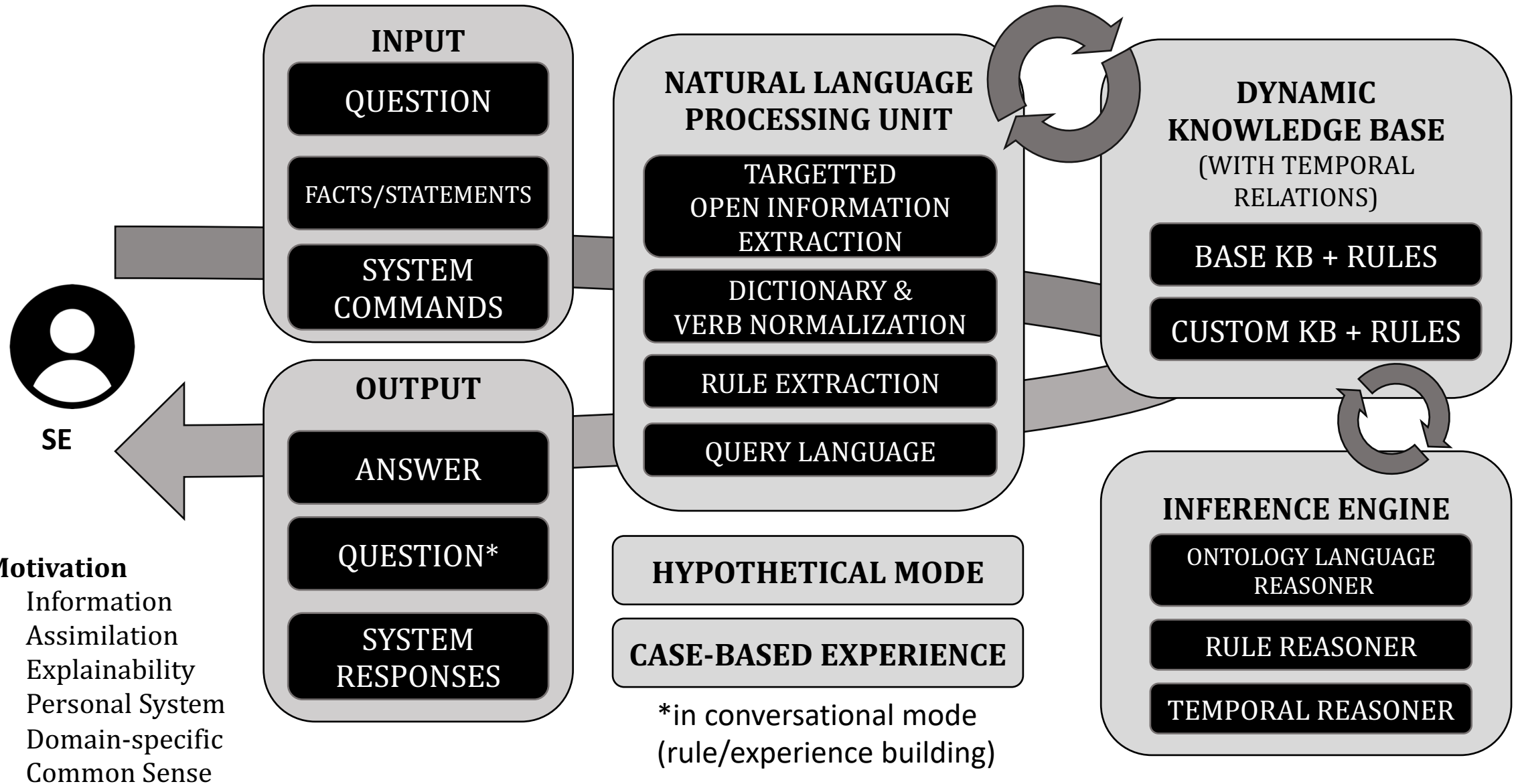
Ben Brumback (George Mason University)

Youngyo Na (Rowan University)

# SEVA: A Systems Engineer's Virtual Assistant



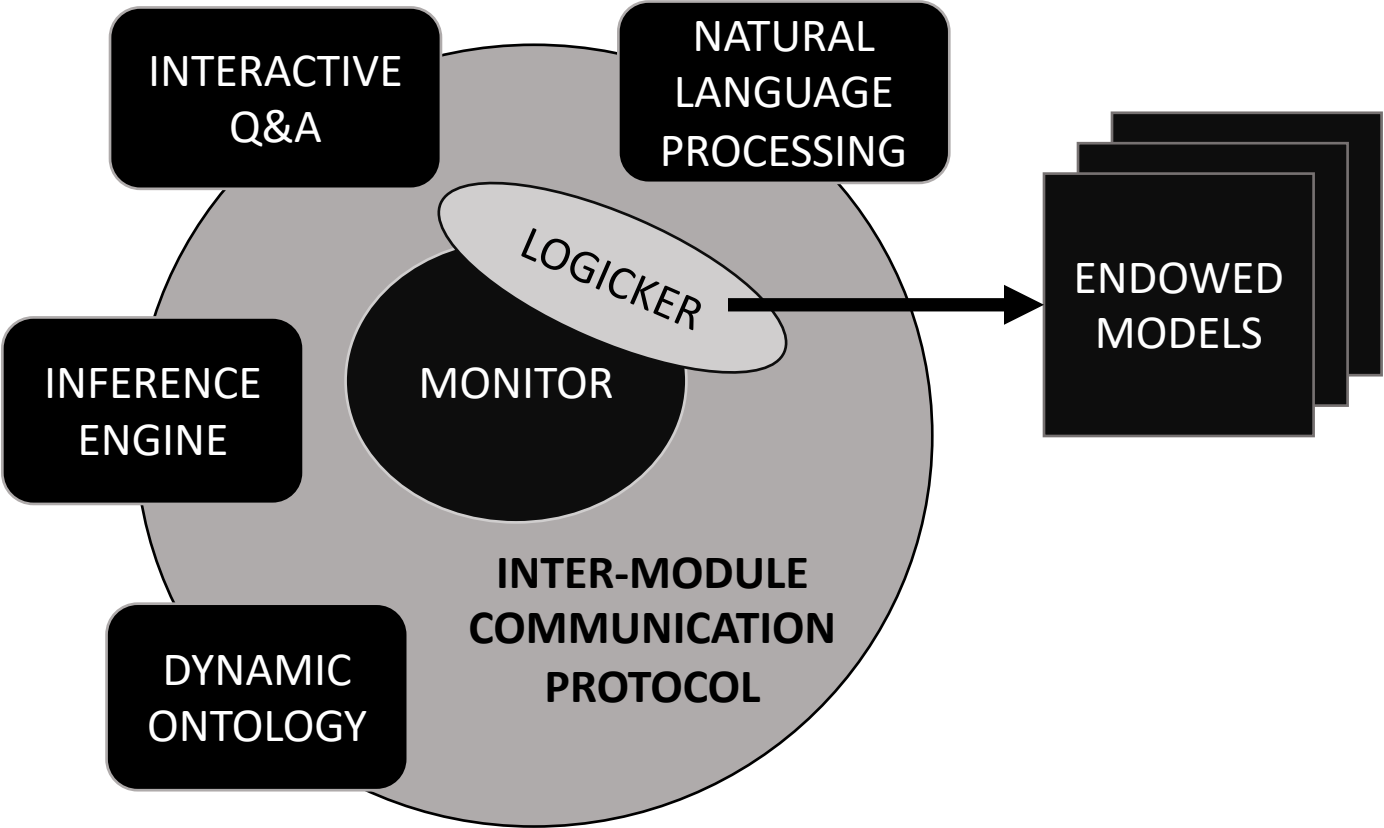
# SEVA: ARCHITECTURAL COMPONENTS



## Motivation

- Information Assimilation
- Explainability
- Personal System
- Domain-specific Common Sense

# SEVA: SYSTEM CONCEPT



# SEVA: Knowledge Representation Example

**Sentence:** Instruments have thermal zones

Subject      Verb/Predicate      Object

NLP

**First Order Logic (FOL) Formalism:**

$\forall x \exists y \text{ Instrument}(x) \wedge \text{Thermal Zone}(y) \Rightarrow \text{have}(x, y)$

Concepts: Instrument, Thermal Zone (like classes in OOP)

Instances: x, y

Relationship: have

KB

$\forall$  = for all  
 $\exists$  = there exists

# SEVA: Knowledge Representation Example

ABox

***Instrument(STI)***

STI is an instrument

***part of(MassSpectrometerAP8717, STI)***

MassSpectrometerAP8717 is a part of STI

Instantiations and relation  
between instances

TBox

***Conduit ≡ Pipe***

Conduit is same as Pipe

***MassSpectrometer ⊆ Spectrometer***

Subclass relationship

Concepts and relations

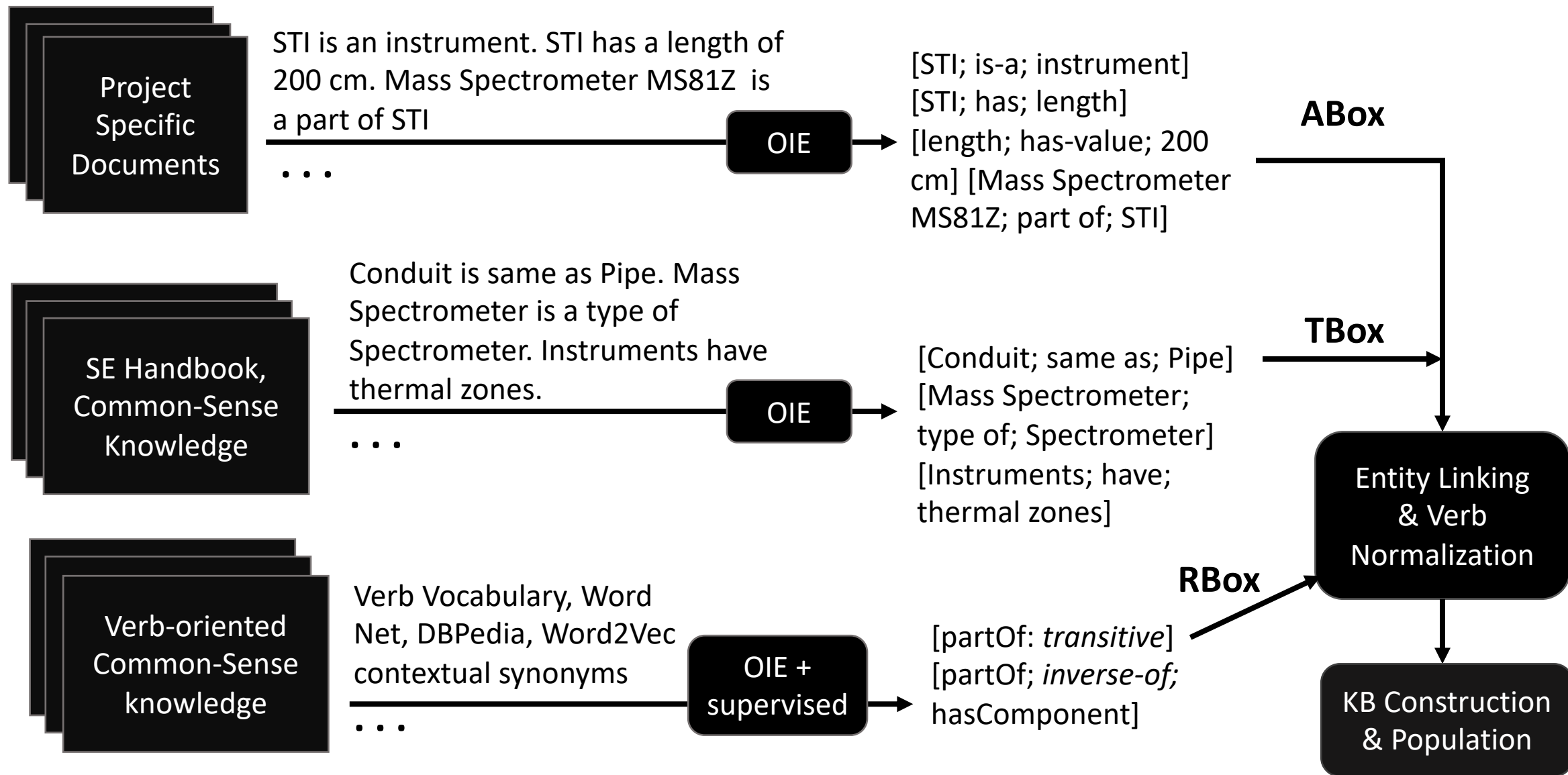
RBox

***partOf ◦ partOf ⊆ partOf***  
Transitive property of the role

Relations between  
relations

***partOf ≡ hasComponent<sup>-</sup>***  
Inverse property of two roles

# SEVA: Types of Knowledge



# Open Information Extraction

- Identifies wide range of domain-independent relations
- Traditional Information Extraction: uses predetermined templates

Examples: Stanford Open IE, Open IE by AI2, ClausIE

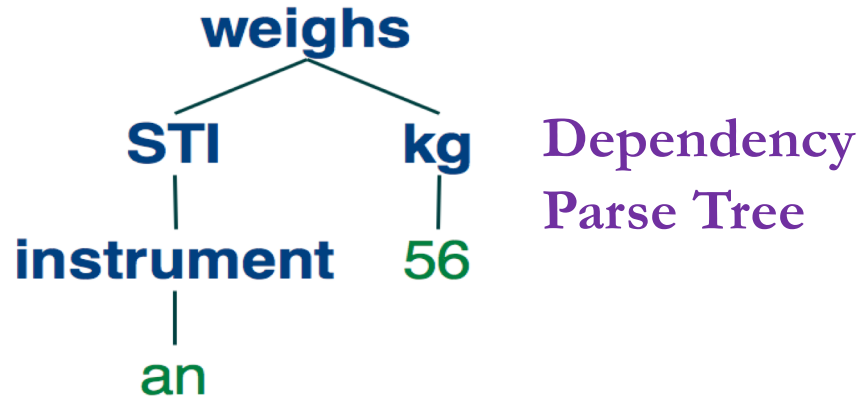
- For the experiment the domain language complexity is reduced – we work with only **simple English** sentences



# NLP Basics

*Sentence:*

**STI**, an instrument, weighs 56 kg

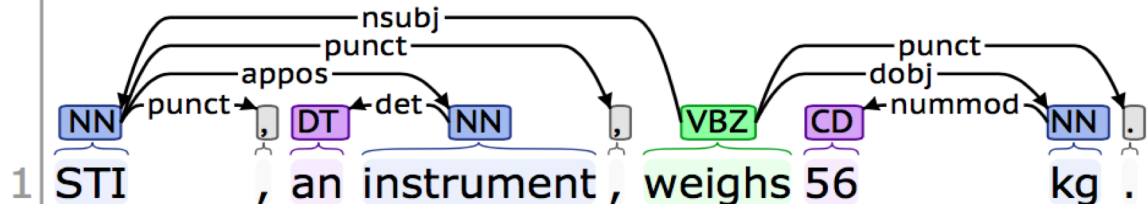


## Part-of-Speech:

NN = Noun  
VBZ = Verb form



## Basic Dependencies:



# SEVA - Targeted Open IE (TOIE)

Pattern Matching on the dependency tree

Extracting relations of type:

**is-a**, **transitive-verb**, **has-property**,  
**has-value**

A set of universal dependencies:

**nsubj**, **dobj**, **case**, **nmod**, **compound**,  
**amod**

Chunking example:

"NP : { (<CD.\*> | <JJ.\*>) <NN.\*>+ } "

NLTK, Stanford CoreNLP, POS Tagger

# SEVA-TOIE: Comparison

Input Sentence	Stanford Open IE	ClausIE	AI2 Open IE	SEVA-TOIE
STI, an instrument, has a 2500 pixel CCD detector	(“STI” “has” “2500 pixel CCD detector”) <i>incomplete/missing information: “STI is an instrument”</i>	(“STI” “is” “an instrument”) (“STI” “has” “a 2500 pixel CCD detector”)	(STI; has; a 2500 pixel CCD detector) (STI; [is]; an instrument)	(STI; has; CCD detector) (STI; is-a; instrument) (CCD detector; has-property; 2500 pixel)
STI is an instrument with a TRL value of 5	(“STI” “is” “instrument”) (“STI” “is instrument with” “TRL value of 5”) (“instrument” “is with” “TRL value of 5”) (“STI” “is instrument with” “TRL value”)	(“STI” “is” “an instrument with a TRL value of 5”) (“STI” “is” “an instrument”)	(STI; is; an instrument with a TRL value of 5) (STI; is an instrument with; a TRL value)	(STI; is; instrument) (instrument; has-property; TRL value) (TRL value; has-value; 5)
STI is scheduled for acoustic testing on July 3, 2015 from 2:00PM to 6:00PM.	<b>produced 13 triples</b> (“STI” “is” “scheduled”) (“STI” “is scheduled for” <b>X</b> ) X = various combinations of remaining sentence	(“STI” “is scheduled” “for acoustic testing on July 3 2015 from 2:00 PM”) (“STI” “is scheduled” “for acoustic testing on July 3 2015 to 6:00 PM”) (“STI” “is scheduled” “for acoustic testing on July 3 2015”)	(STI; is scheduled; for acoustic testing) <i>incomplete/missing information: “on July 3, 2015 from 2:00PM to 6:00PM”</i>	(STI; is; scheduled) (scheduled; for; testing) (testing; has-property; acoustic) (testing; has-value; 3 2015 July) (scheduled; from; 2:00 PM) (scheduled; to; 6:00 PM)]

Table 1: Output from Open Information Extractors. **Red** colored extractions are incomplete, incorrect, or noisy; **Blue** colored extractions need to be further granularized for ontology population.

**THANK YOU!**