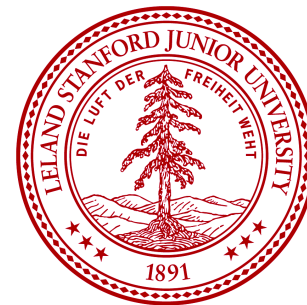


A computational method for the extraction of pharmacogenomic relationships from text

Adrien Coulet^{1,2}, Nigam Shah², Yael Garten²,
Mark Musen², Russ Altman²

1 LORIA, INRIA Nancy – Grand-Est

2 Stanford University



The NCBO and PharmGKB

- A joint project

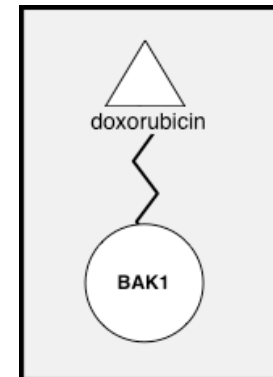


- Content of PharmGKB

- Current:

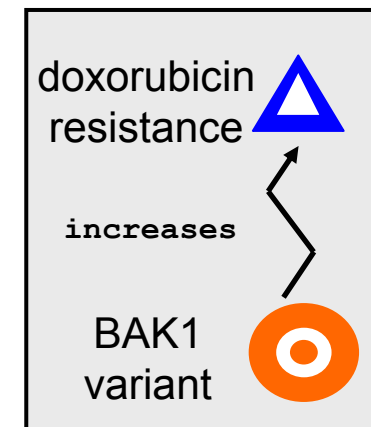
pharmacogenomics (PGx) relationships

Gene – Drug ; Gene – Disease ; Drug – Disease



- Goal:

to provide more precise relationships



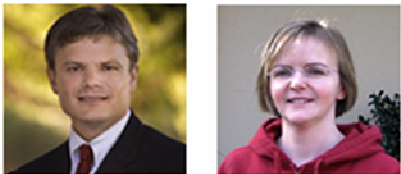
Population of PharmGKB

Sentence 1 : Doxorubicin induces BAK1 activity.

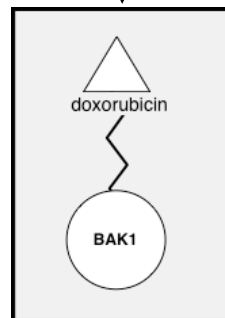
Sentence 2 : Doxorubicin transcriptionally activates BAK1.

Sentence 3 : BAK1 gene polymorphism affects doxorubicin resistance.

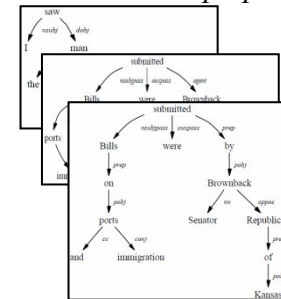
Scientific literature



PharmGKB curators

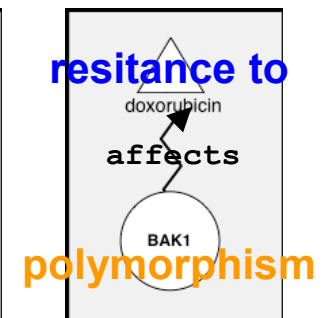
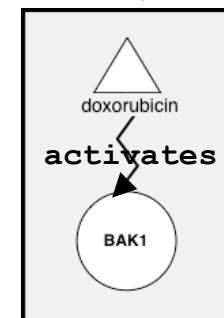
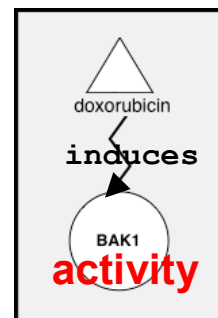


Dependency
Graph parsing



Dependency Graphs
of sentences

Relation
extraction



Population of PharmGKB

Sentence 1 : Doxorubicin induces BAK1 activity.

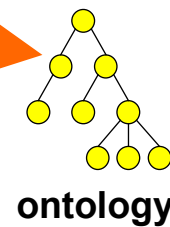
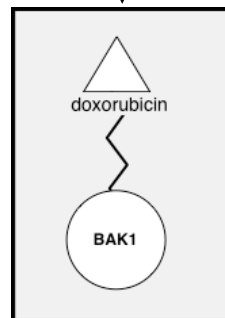
Sentence 2 : Doxorubicin transcriptionally activates BAK1.

Sentence 3 : BAK1 gene polymorphism affects doxorubicin resistance.

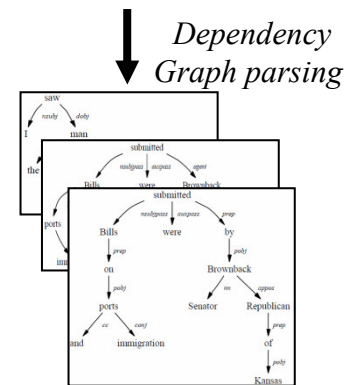
Scientific literature



PharmGKB curators

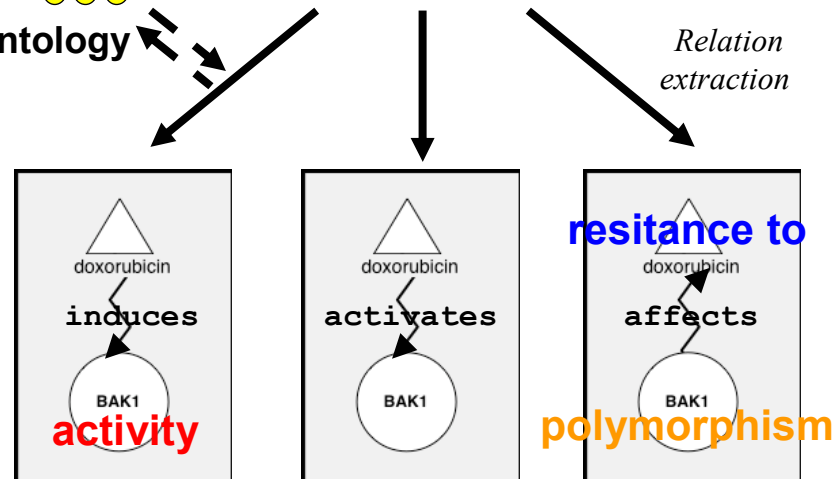


ontology



Dependency Graph parsing

Dependency Graphs of sentences



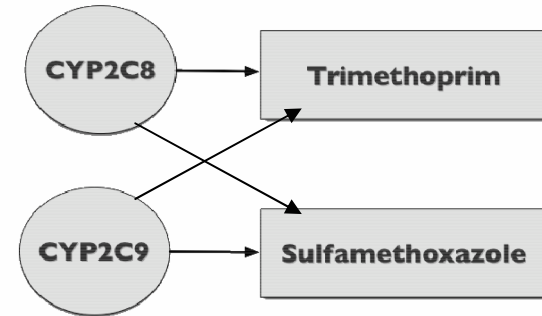
Outline

1. Limitations of co-occurrences
2. Construction of semantic network
 1. Algorithm to extract raw relationships
 2. Semi-automated ontology building
 3. Comprehensive knowledge network from 1 & 2

Limitations of co-occurrence (that we wanted to solve)

1. Avoid false positive connections

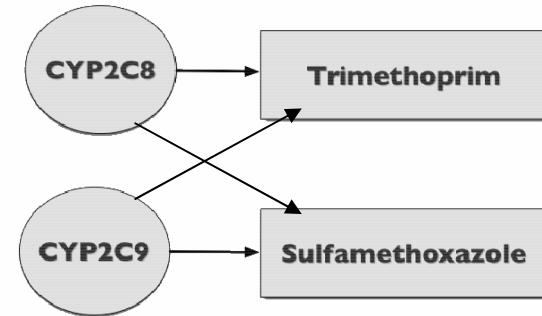
“Trimethoprim inhibits activity of CYP2C8 while sulfamethoxazole inhibits CYP2C9 activity.”



Limitations of co-occurrence (that we wanted to solve)

1. Avoid false positive connections

“Trimethoprim inhibits activity of CYP2C8 while sulfamethoxazole inhibits CYP2C9 activity.”



2. Characterize fine-grain semantics of relationships

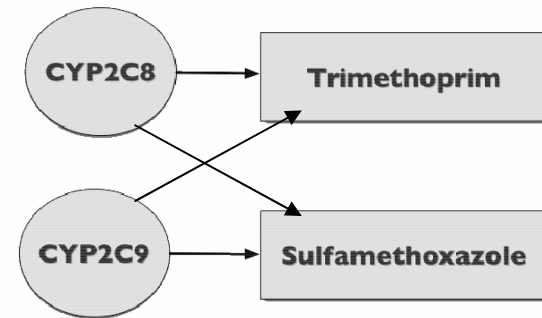
“CYP3A4 mRNA expression was increased significantly by rifampicin exposure in human hepatocytes.”



Limitations of co-occurrence (that we wanted to solve)

1. Avoid false positive connections

“Trimethoprim inhibits activity of CYP2C8 while sulfamethoxazole inhibits CYP2C9 activity.”



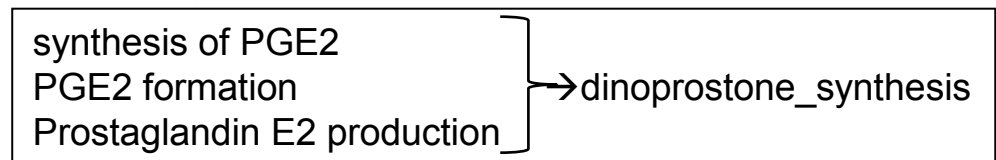
2. Characterize fine-grain semantics of relationships

“CYP3A4 mRNA expression was increased significantly by rifampicin exposure in human hepatocytes.”

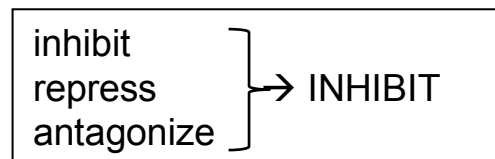


3. To consolidate synonyms (normalize):

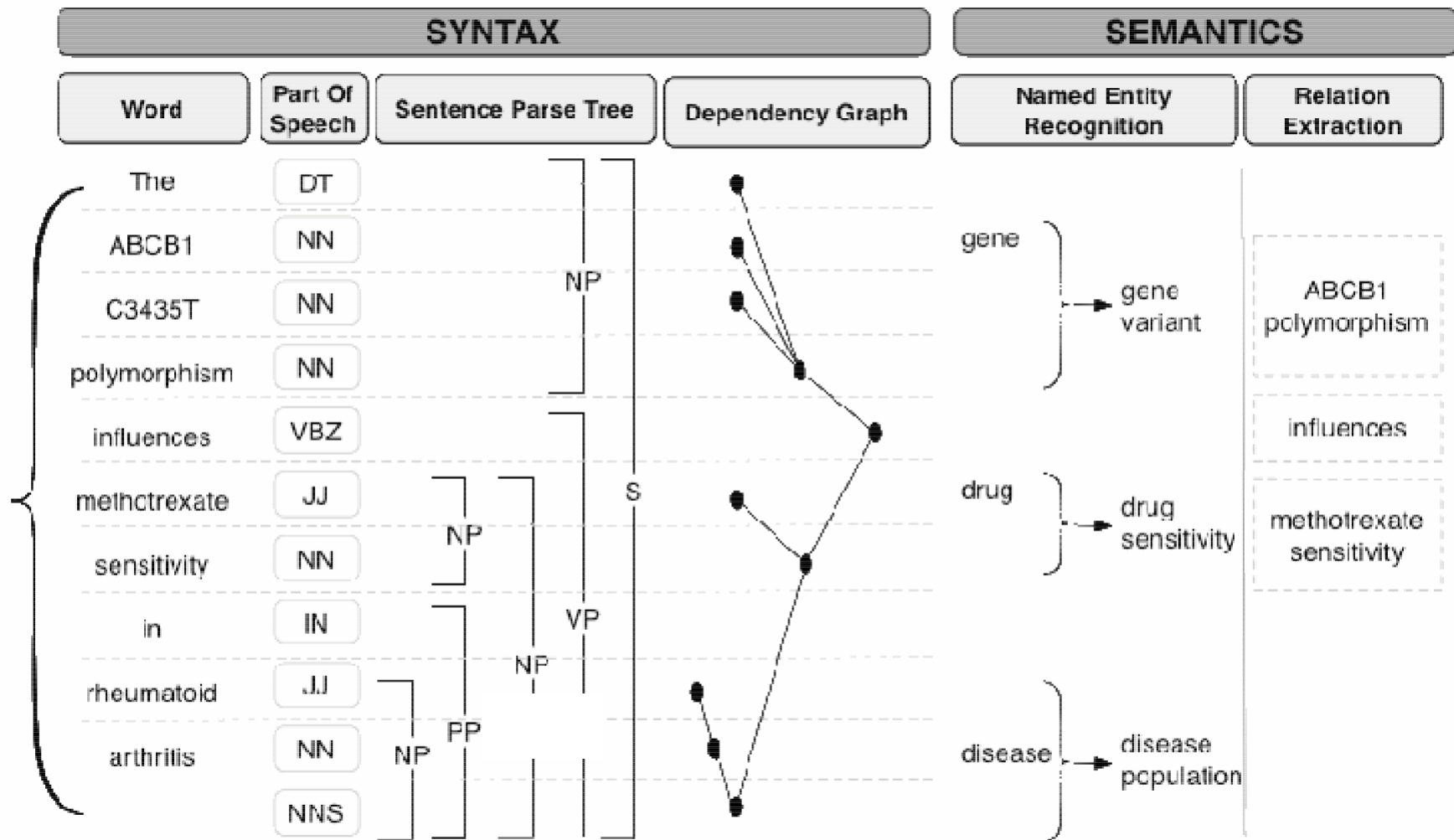
- Between complex entity names:



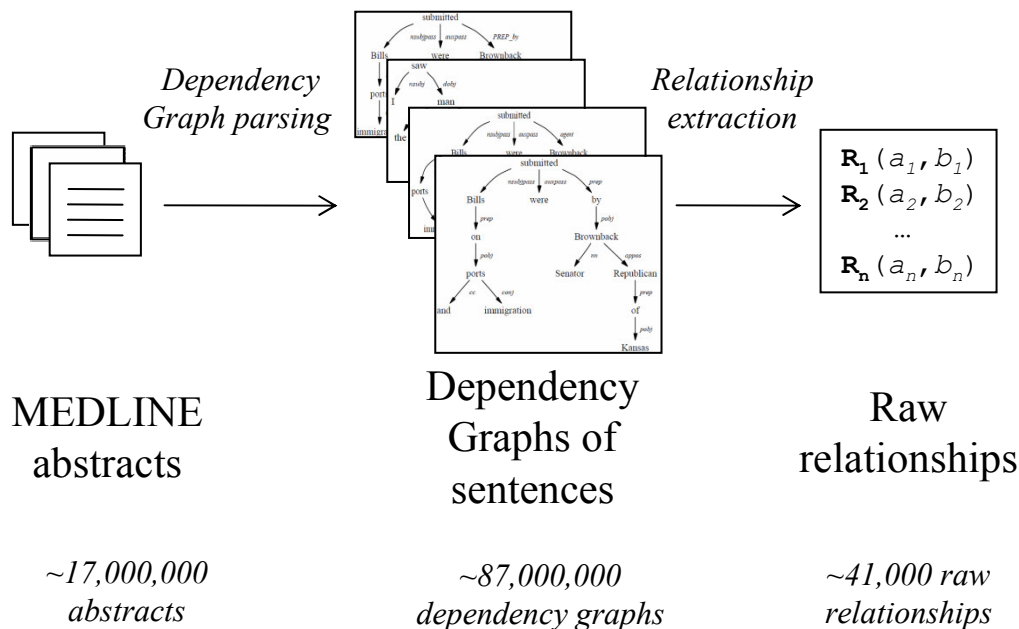
- Between relationships:



Several steps of text processing enable extracting relationship semantics



The method extracts high quality typed relationships



Evaluation:

Randomly selected 220 raw relationships:
classified into 3

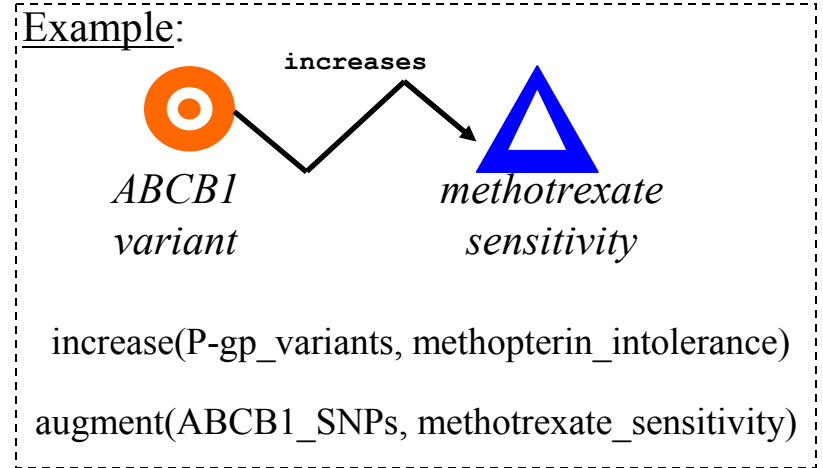
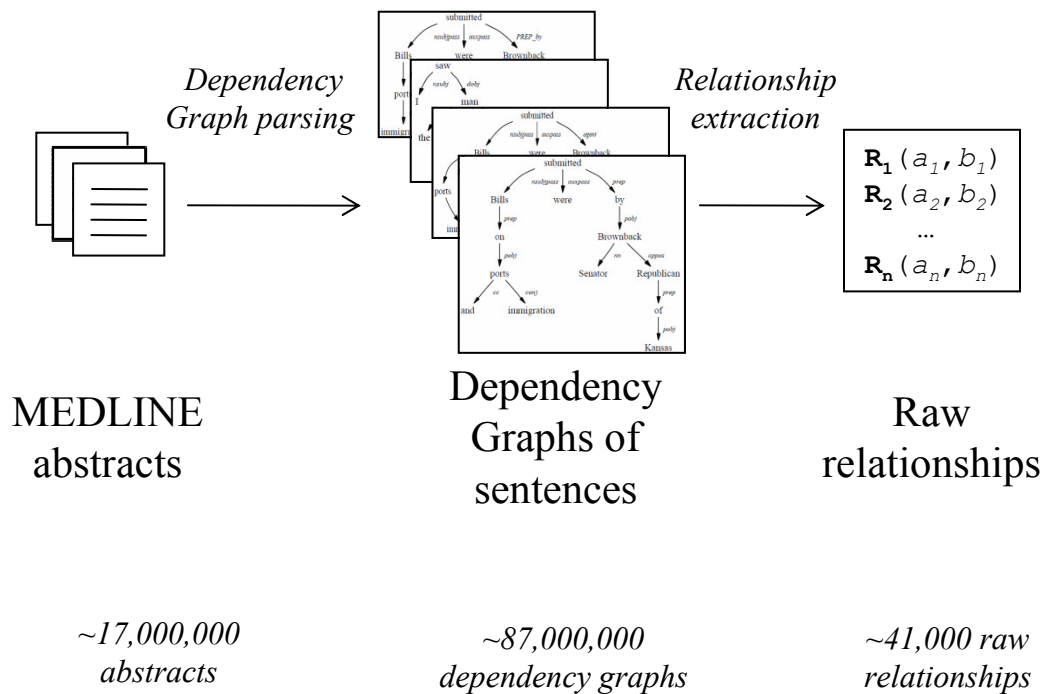
“polymorphisms in VKORC1 are associated with warfarin dose.”

- associated(VKORC1_polymorphisms, warfarin_dose)
= true and complete
- associated (VKORC1_polymorphisms, warfarin)
= true and incomplete
- polymorphisms (VKORC1, warfarin_dose)
= false

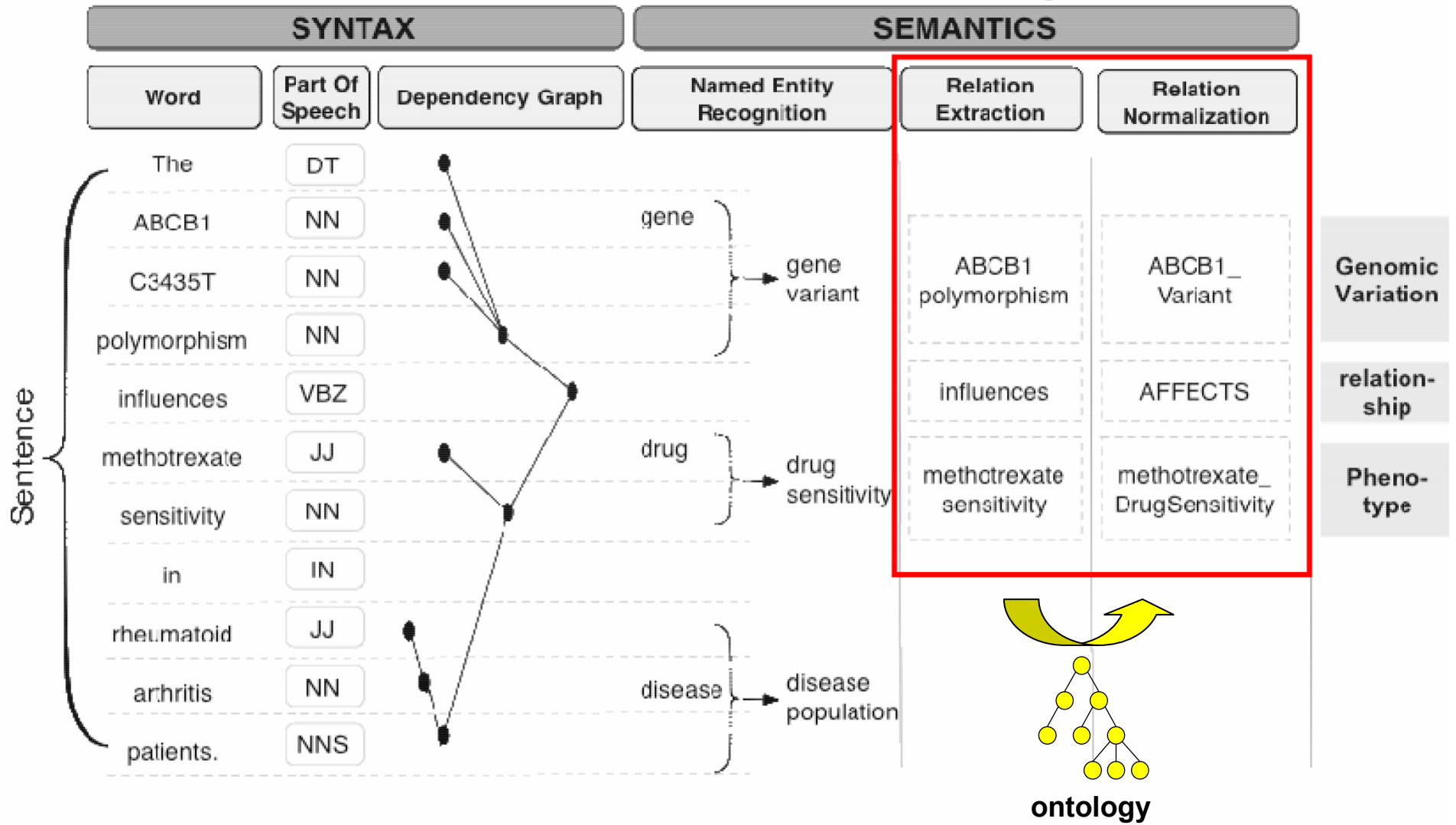
Results:

- 87.7% were complete or incomplete true positives
 - 70% true and complete
 - 17.7% true and incomplete
- 12.3% were false positives

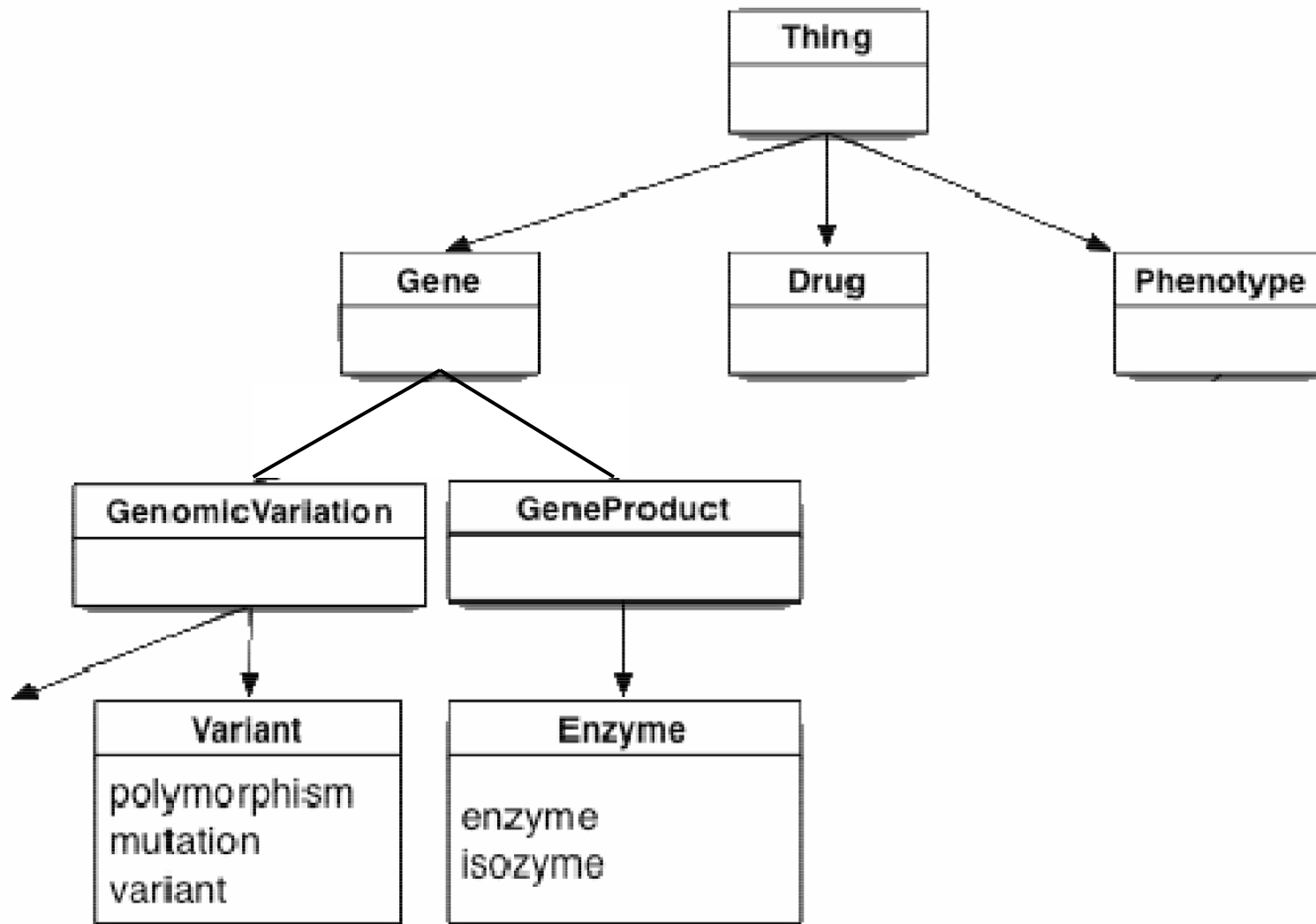
Issue: we extracted heterogeneous relationships



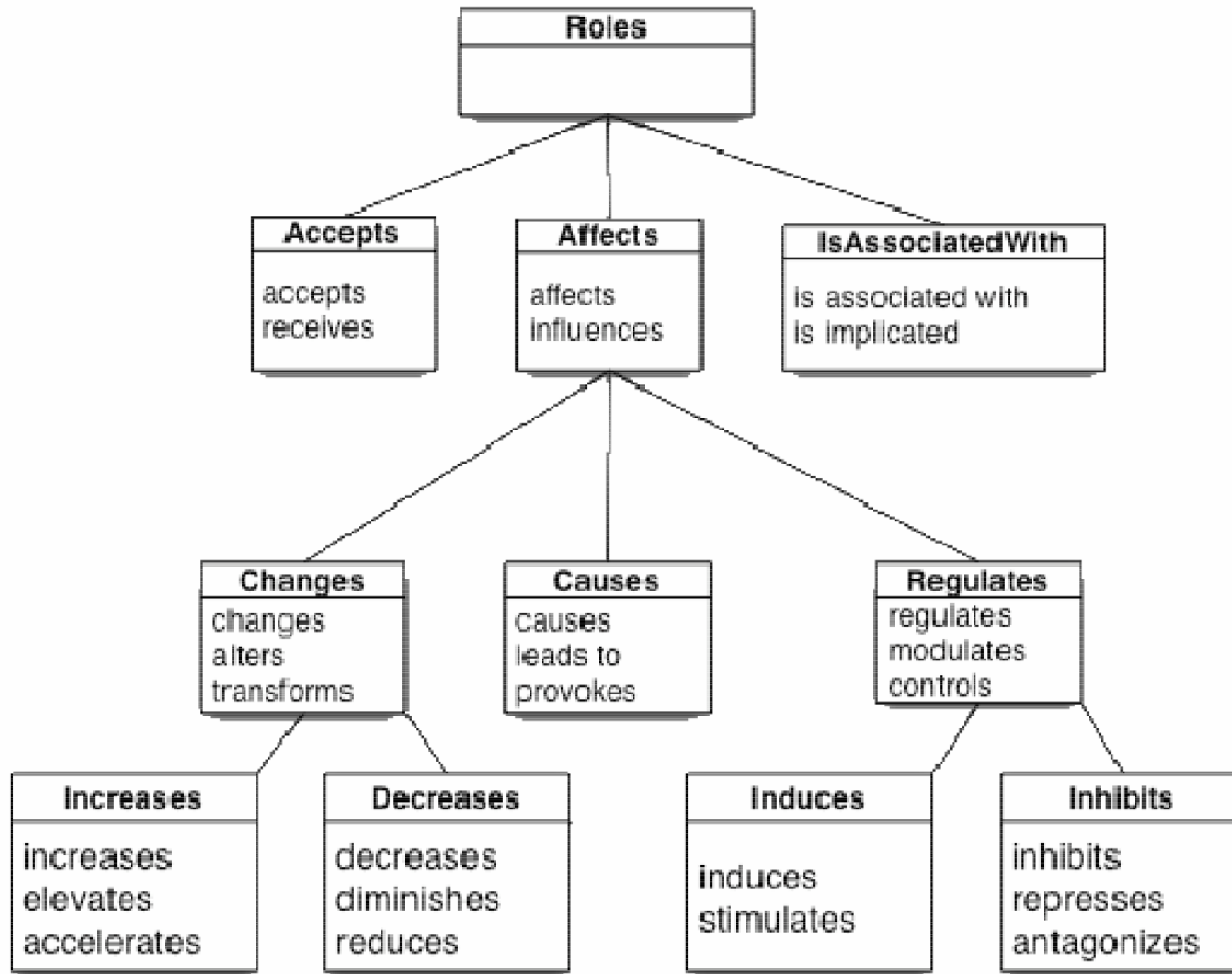
We built and use an ontology to normalize relationships



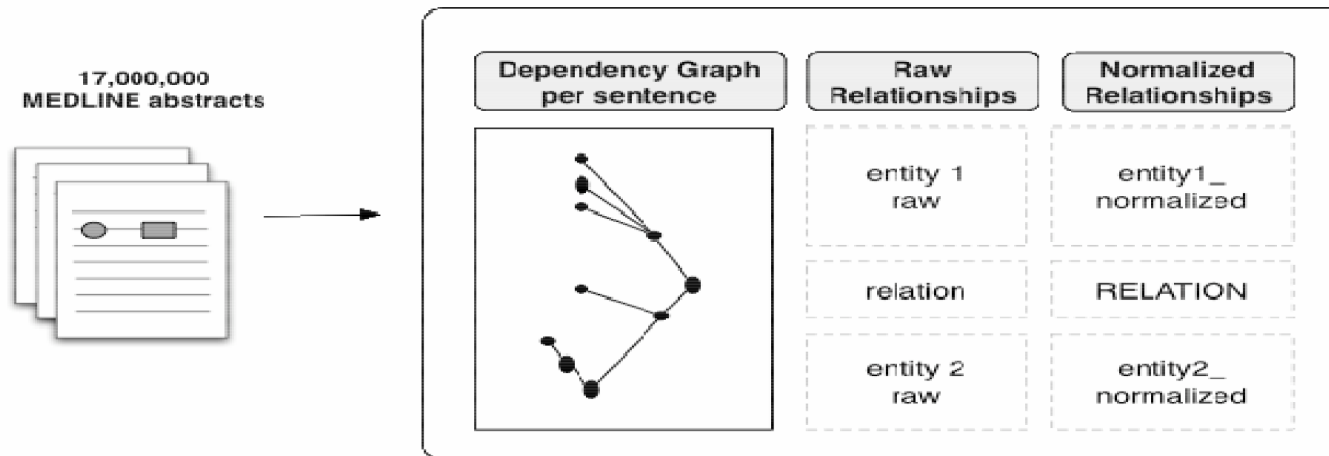
An ontology organizes the “world” into concepts and roles (1/2)



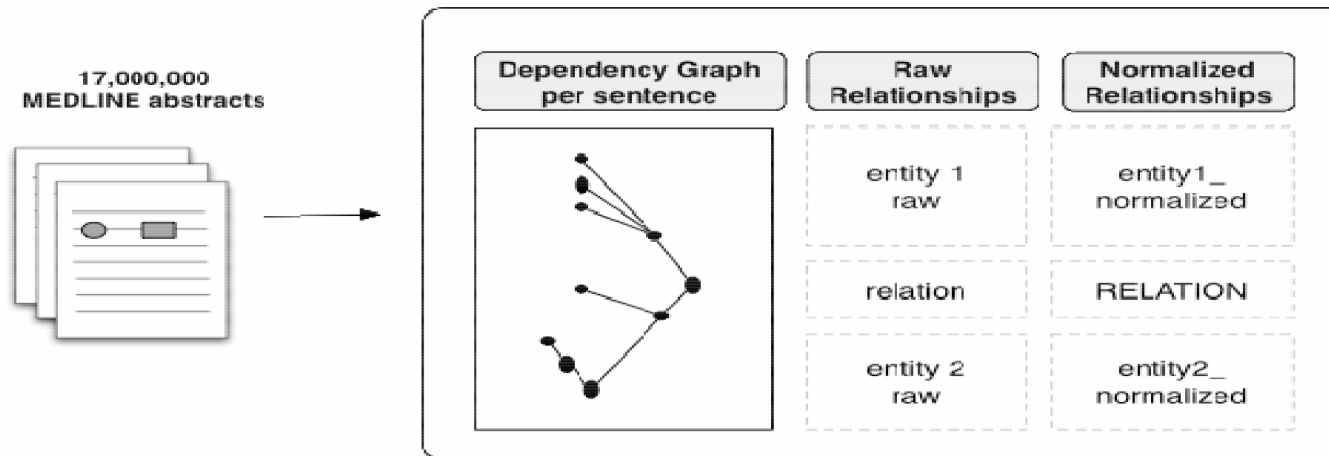
An ontology organizes the “world” into concepts and roles (2/2)



We manually created a PGx ontology “bottom-up”

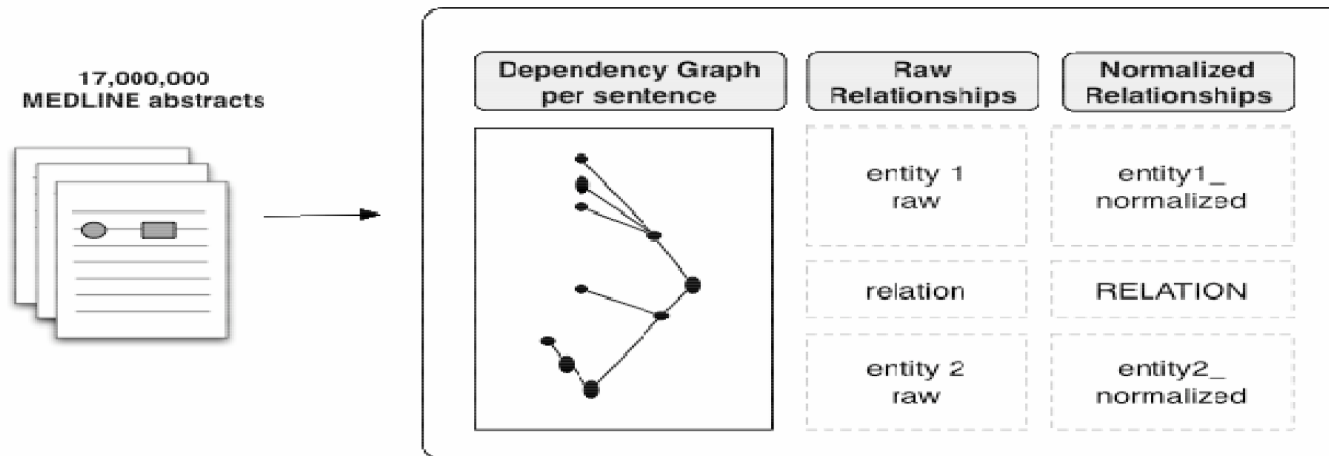


We manually created a PGx ontology “bottom-up”



Relationship types	Entities modified by		
	Genes	Drugs	Phenotypes
2538 associate	1237 <i>gene</i>	377 <i>metabolism</i>	304 <i>cell</i>
1017 increase	1000 <i>inhibitor</i>	358 <i>activity</i>	114 <i>line</i>
985 inhibit	935 <i>polymorphism</i>	298 <i>inhibitor</i>	101 <i>patient</i>
825 induce	775 <i>expression</i>	267 <i>effect</i>	71 <i>risk</i>
763 metabolize	773 <i>activity</i>	263 <i>administration</i>	35 <i>tissue</i>
666 involve	689 <i>mutation</i>	246 <i>channel</i>	34 <i>specimen</i>
643 reduce	685 <i>genotype</i>	242 <i>treatment</i>	33 <i>case</i>
547 catalyze	393 <i>inhibition</i>	193 <i>antagonist</i>	27 <i>treatment</i>
515 cause	329 <i>level</i>	178 <i>concentration</i>	26 <i>rate</i>
509 affect	245 <i>gene_mutation</i>	172 <i>dose</i>	26 <i>effect</i>

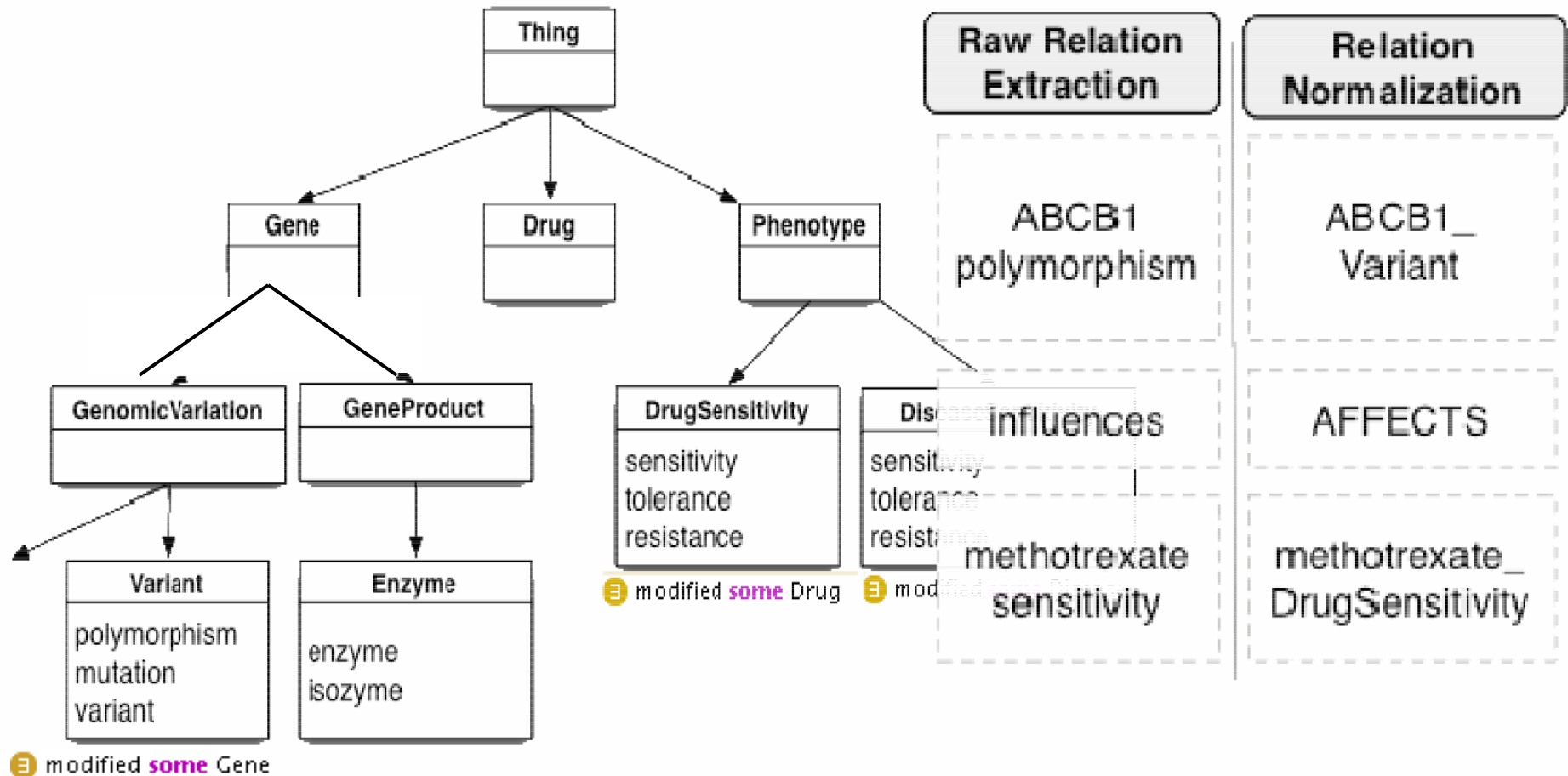
We manually created a PGx ontology “bottom-up”



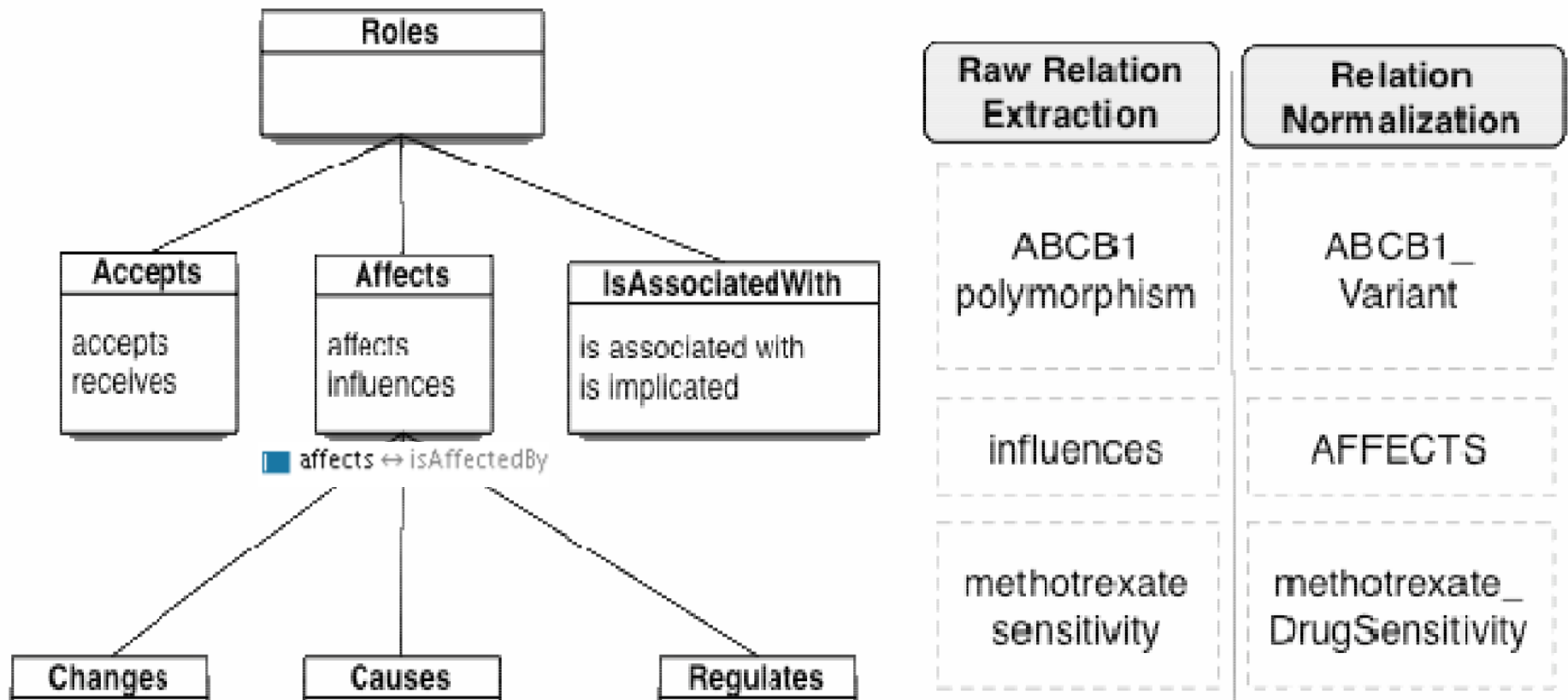
Relationship types	Entities modified by		
	Genes	Drugs	Phenotypes
2538 associate	1237 <i>gene</i>	377 <i>metabolism</i>	304 <i>cell</i>
1017 increase	1000 <i>inhibitor</i>	358 <i>activity</i>	114 <i>line</i>
985 inhibit	935 <i>polymorphism</i>	298 <i>inhibitor</i>	101 <i>patient</i>
825 induce	775 <i>expression</i>	267 <i>effect</i>	71 <i>risk</i>
763 metabolize	773 <i>activity</i>	263 <i>administration</i>	35 <i>tissue</i>
666 involve	689 <i>mutation</i>	246 <i>channel</i>	34 <i>specimen</i>
643 reduce	685 <i>genotype</i>	242 <i>treatment</i>	33 <i>case</i>
547 catalyze	393 <i>inhibition</i>	193 <i>antagonist</i>	27 <i>treatment</i>
515 cause	329 <i>level</i>	178 <i>concentration</i>	26 <i>rate</i>
509 affect	245 <i>gene_mutation</i>	172 <i>dose</i>	26 <i>effect</i>

237 concepts
76 roles

We use the ontology to normalize the raw relationship (subject, relation and object)



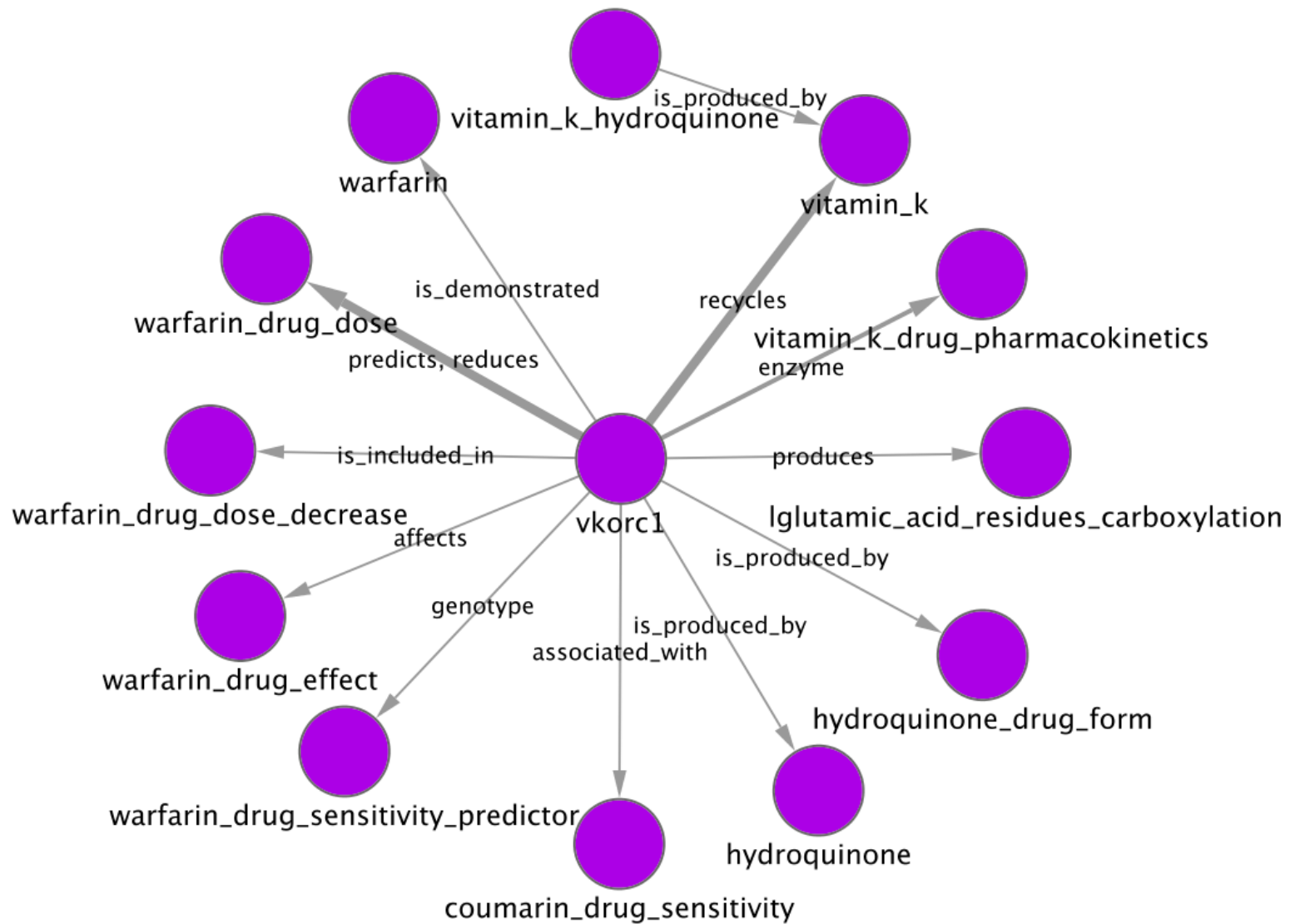
We use the ontology to normalize the raw relationship (subject, relation and object)



Example: two sentences but one fact

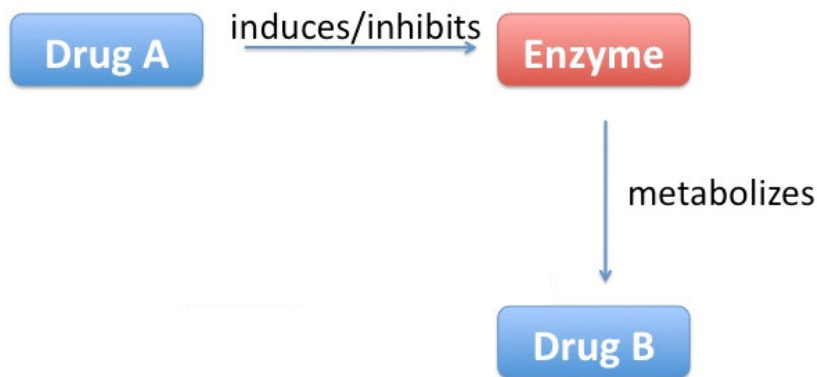
raw text	sentence	The ABCB1 C3435T polymorphism influences methotrexate sensitivity in rheumatoid arthritis patients.	A variant C3435T allele of the MDR1 gene affects methotrexate tolerability.
	raw relationship	entity 1	ABCB1 polymorphism
	relationship	influences	affects
	entity2	methotrexate sensitivity	methotrexate tolerability
normalized relationship	entity 1	ABCB1_Variant	ABCB1_Variant
	relationship	AFFECTS	AFFECTS
	entity2	methotrexate_DrugSensitivity	methotrexate_DrugSensitivity

Example of network (1/3): VKORC1



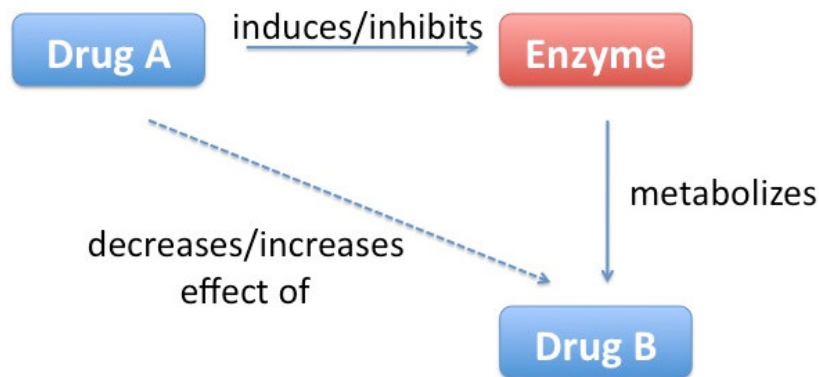
Conclusion

- A method to build semantic network
- Used in PGx:
 - For curation and knowledge summarization
@PharmGKB
 - For knowledge discovery
e.g. Predicting Drug-Drug interaction
=> *Yael Garten's PhD thesis*



Conclusion

- A method to build semantic network
- Used in PGx:
 - For curation and knowledge summarization
@PharmGKB
 - For knowledge discovery
e.g. Predicting Drug-Drug interaction
=> *Yael Garten's PhD thesis*



Questions?

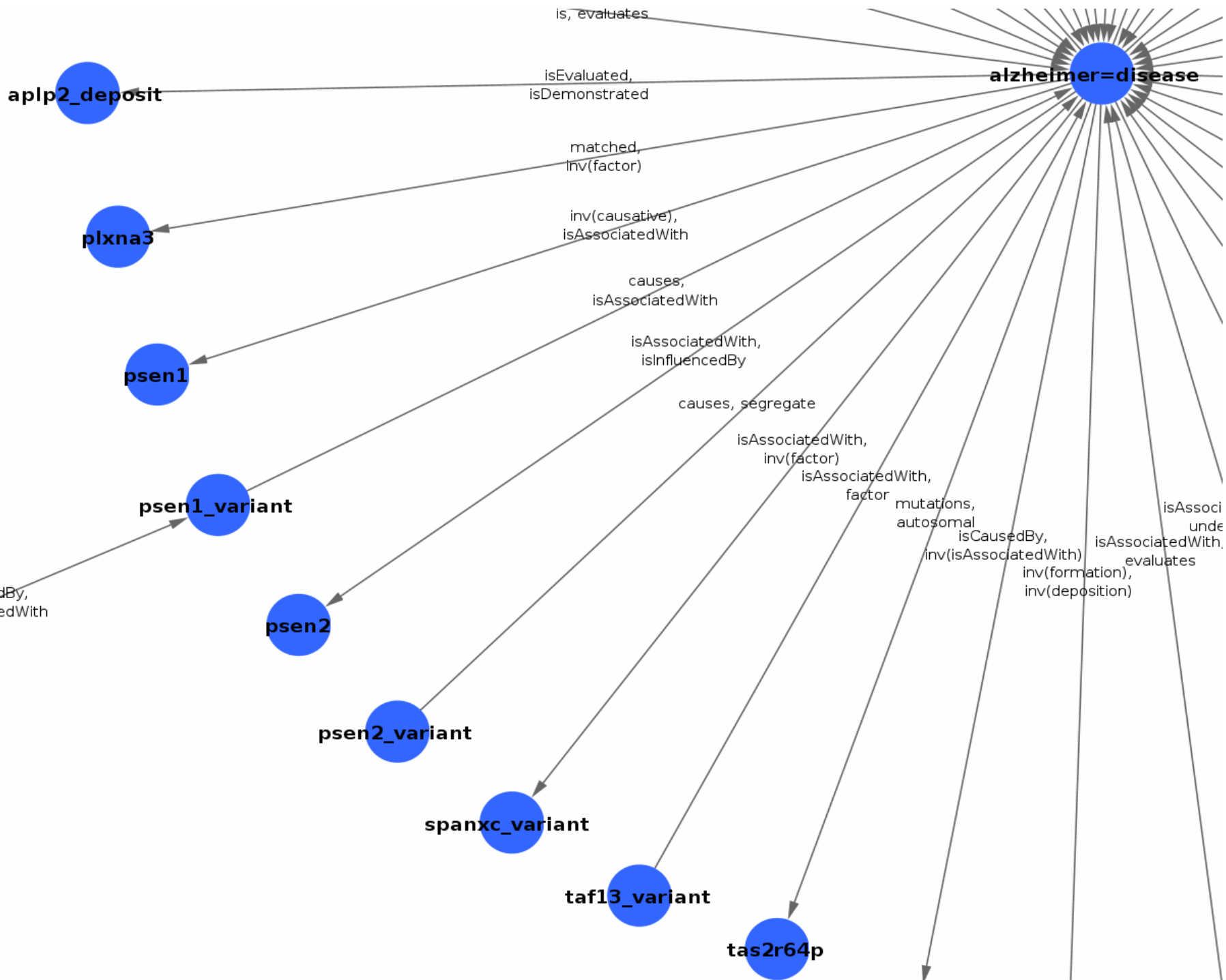
Coulet *et al.* *Journal of Biomedical Informatics*, In Press, 2010

or

adrien.coulet@loria.fr

Thanks

And thanks to Yael Garten for many slides



relationship type

subject

object

inhibited (*VKORC1_expression*, *warfarin*)

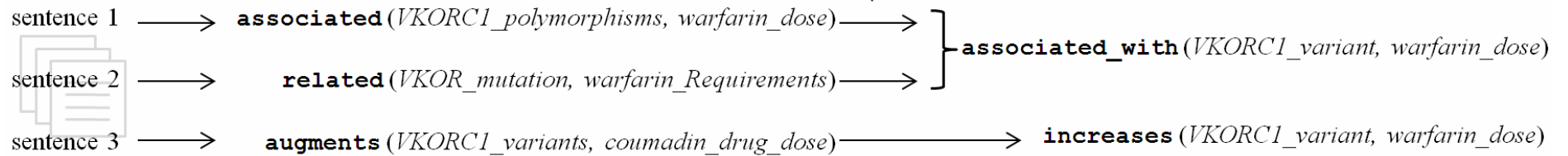
PGx key entity

modified entity

single PGx
key entity

PHARE Ontology

Concepts	Variant hasLabel {variant, polymorphism, mutation} DrugDose hasLabel {dose, requirement}
roles	associated_with hasLabel {associated, related} increases hasLabel {induce, increase}
<i>individuals</i>	<i>VKORC1</i> hasLabel {VKORC1, VKOR} <i>warfarin</i> hasLabel {warfarin, coumadin}



Sentences

Raw
relationships

Normalized
relationships