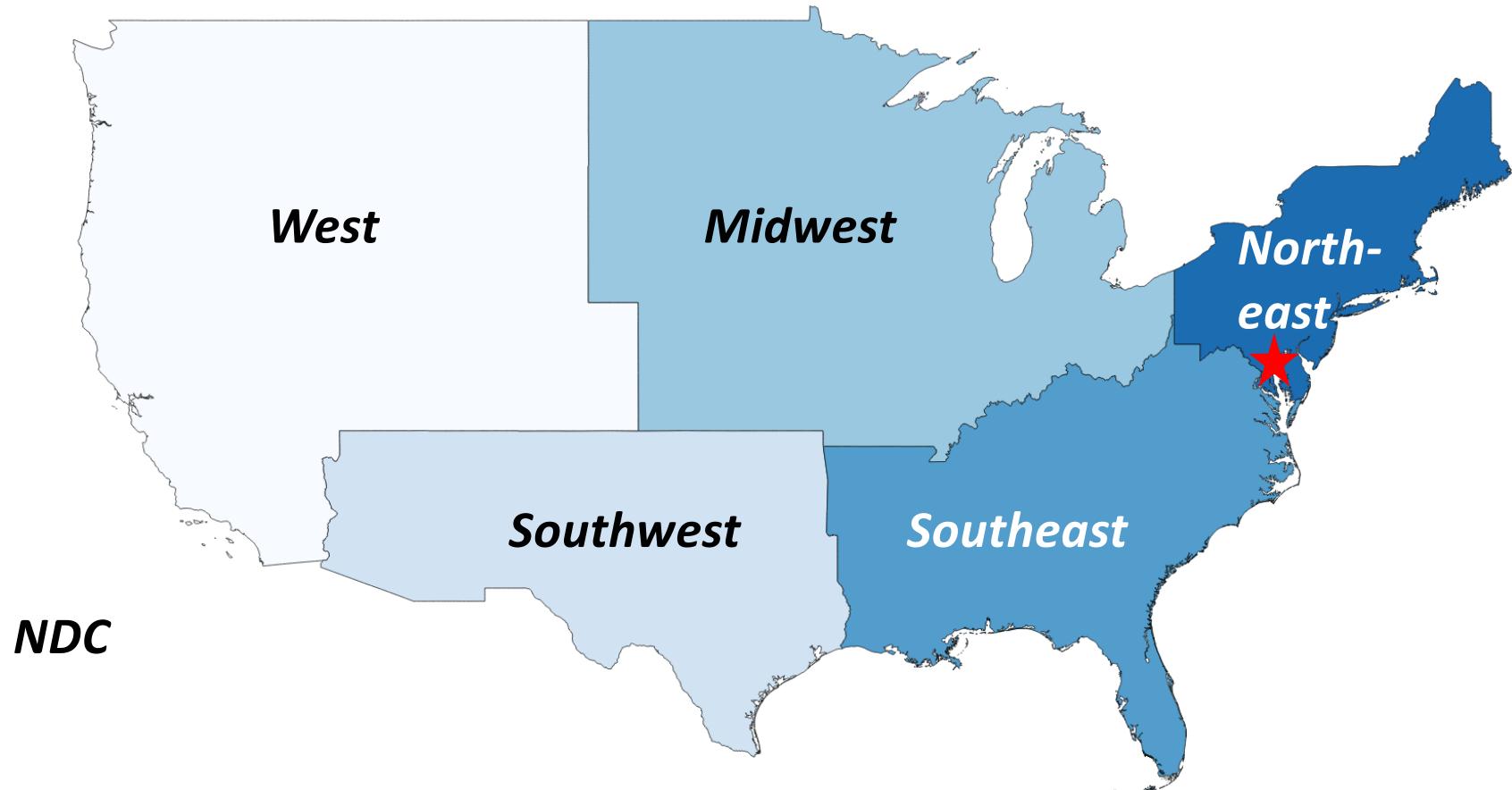


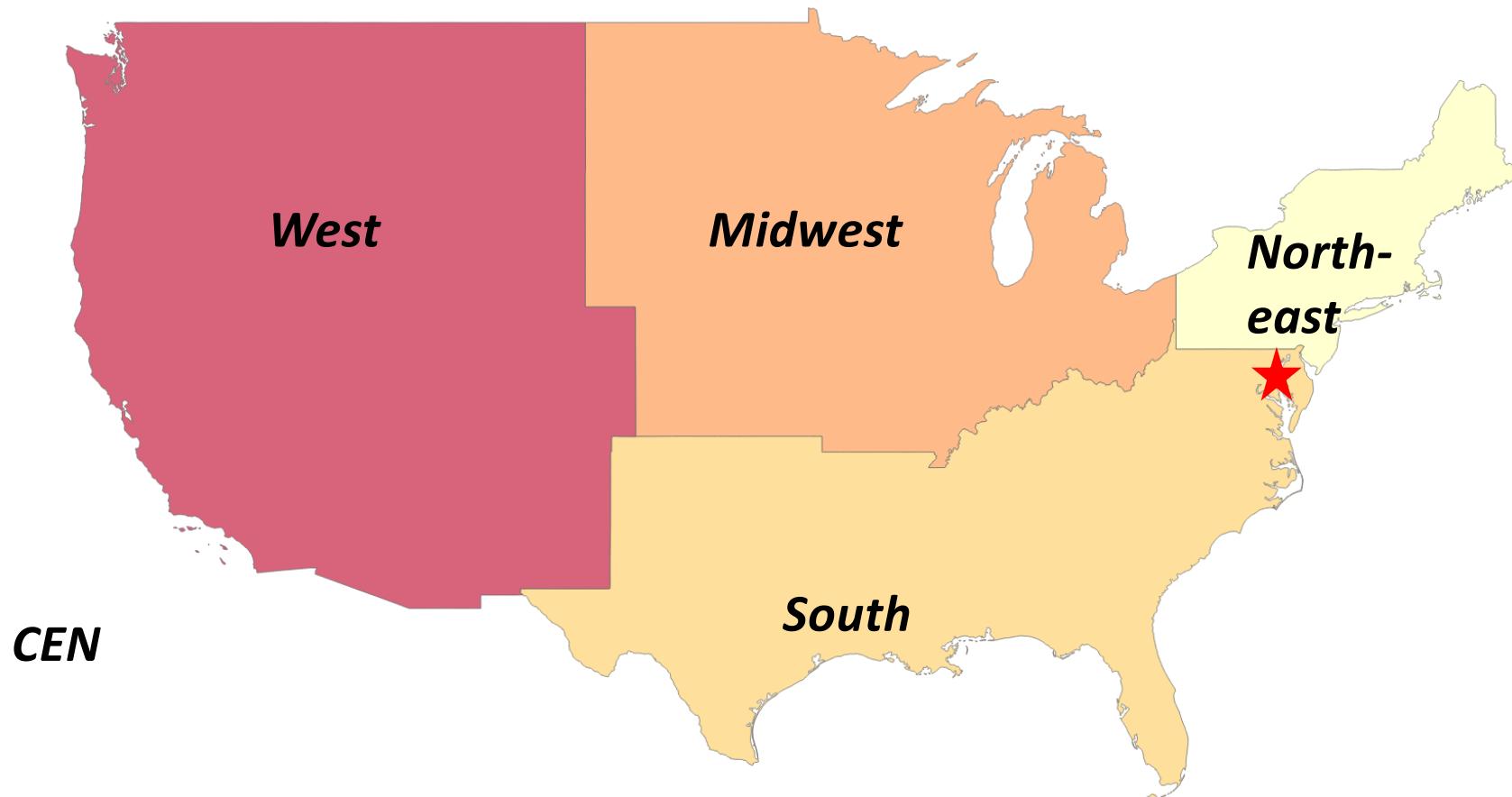
Theory & Practice of Data Cleaning

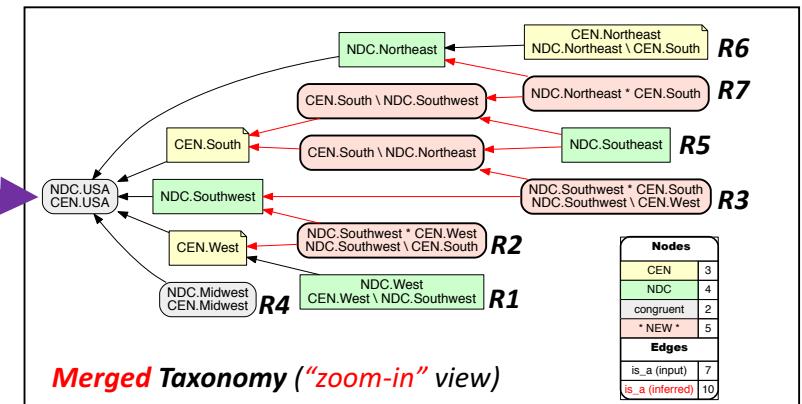
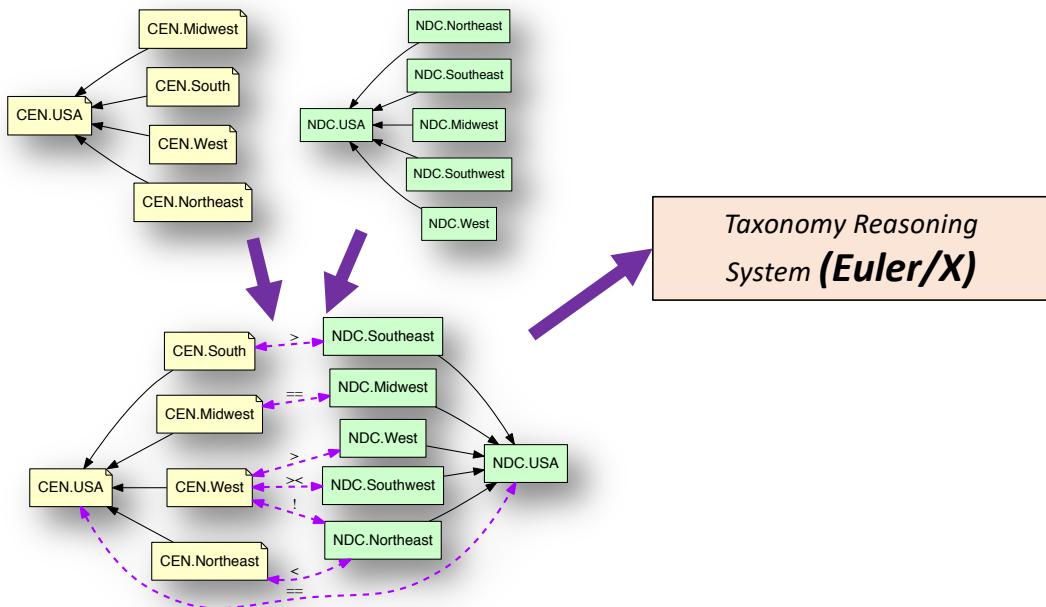
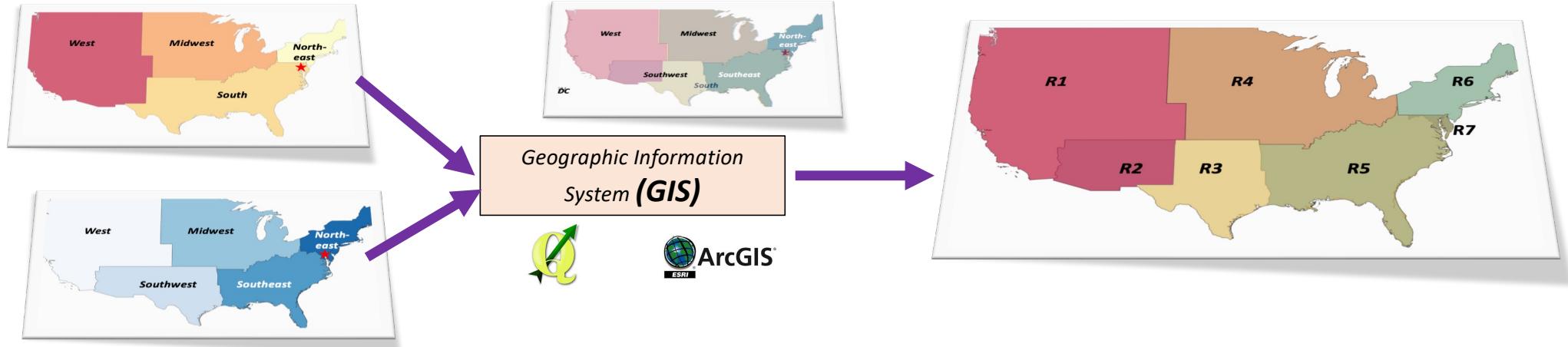
From Controlled Vocabularies to Taxonomy Alignment

Controlled Vocabularies & Taxonomies

- **Tina:** Hey, any recommended signature dish from where you live?
- **Amy:** Oh, definitely the half-smokes from the Northeast! They are these tasty half-pork and half-beef sausages.
- **Tina:** What a coincidence! We have half-smokes in the South, too! Where do you live in the Northeast? New York? Boston?
- **Amy:** Wrong guesses! Where do you live in the South?
- **Tina and Amy together:** Washington, D.C. [The two of them look at each other, confused.]







Overview

- The **Taxonomy Alignment Problem**
 $T_1 + T_2 + A \Rightarrow T_3$ (*ambiguous .. unique .. inconsistent*)
- Model-based **Diagnosis** [Reiter'87]
 - Black-box
- Hybrid Approach
 - Black-box & White-box combined

Optional Readings:

Cheng, Yi-Yun, Nico Franz, Jodi Schneider, Shizhuo Yu, Thomas Rodenhausen, and Bertram Ludäscher. “**Agreeing to Disagree: Reconciling Conflicting Taxonomic Views Using a Logic-Based Approach.**” *Proceedings of the Association for Information Science and Technology* 54, no. 1 (January 1, 2017): 46–56. <https://doi.org/10.1002/pra2.2017.14505401006>.

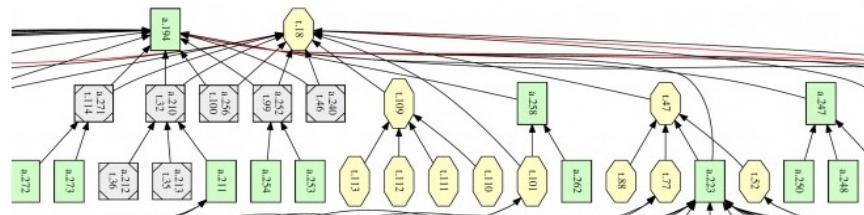
Chen, Mingmin, Shizhuo Yu, Nico Franz, Shawn Bowers, and Bertram Ludäscher. “**A Hybrid Diagnosis Approach Combining Black-Box and White-Box Reasoning.**” In *Rules on the Web. From Theory to Applications*, edited by Antonis Bikakis, Paul Fodor, and Dumitru Roman, 127–41. Lecture Notes in Computer Science. Cham: Springer International Publishing, 2014. https://doi.org/10.1007/978-3-319-09870-8_9.

Meet Prof. Nico Franz: Curator of Insects @ ASU



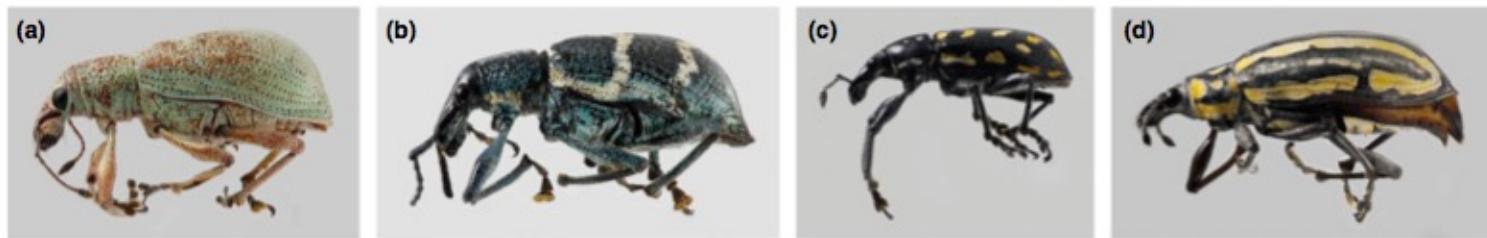
296

N. M. Franz / Cladistics 30 (2014) 294–321



Initial alignment of two weevil classifications with the Euler toolkit

Linked here is an interim result of my attempt to align two influential weevil classifications by Thompson (1992) and Alonso-Zarazaga & Lyal

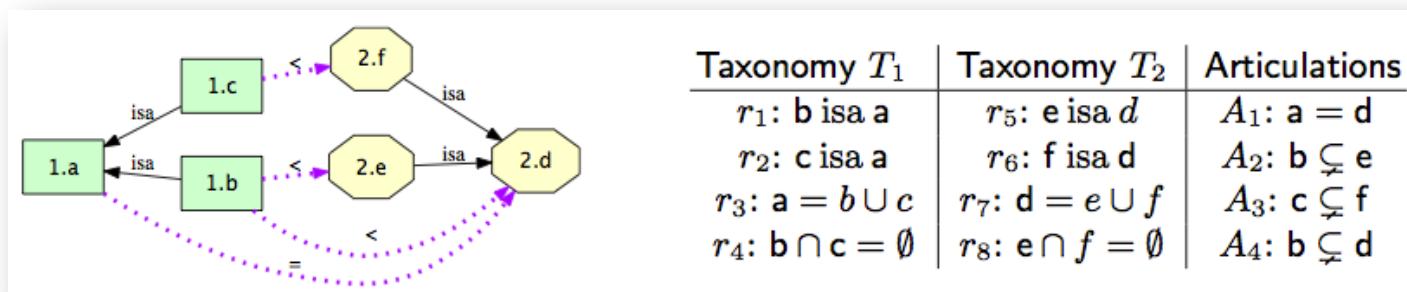


What Nico et al. do for a living ...



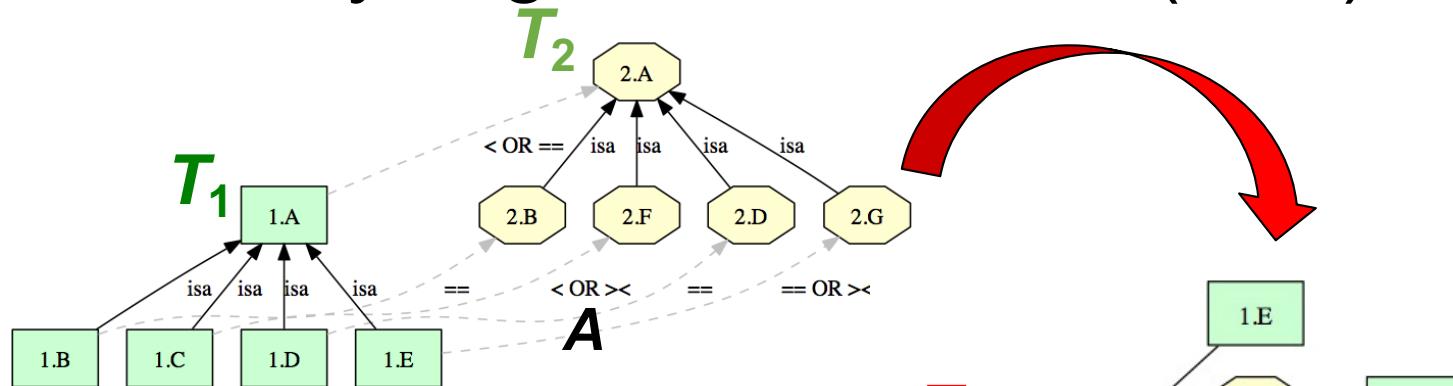
What Nico does for a living (cont'd): The Indoors Part

- First: go fun places, find new bugs, study them ...
 - “Bugs-R-Us” (see [taxonbytes.org](#))
- Now: Compare, **align** and **revise** taxonomies, based on careful observation, “character” data, expertise ...
- Formally:
 - Input: T_1 + T_2 (*taxonomies*) + A (*expert articulations*)

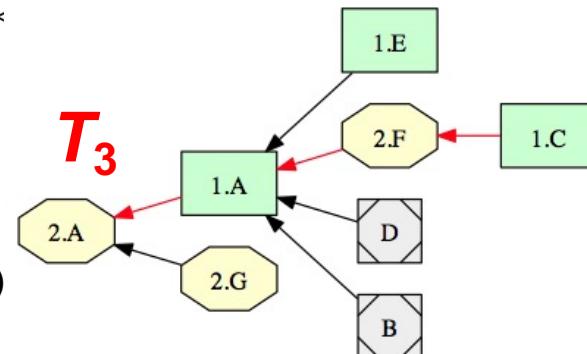


- Output: revised, “merged” taxonomy (-ies) T_3

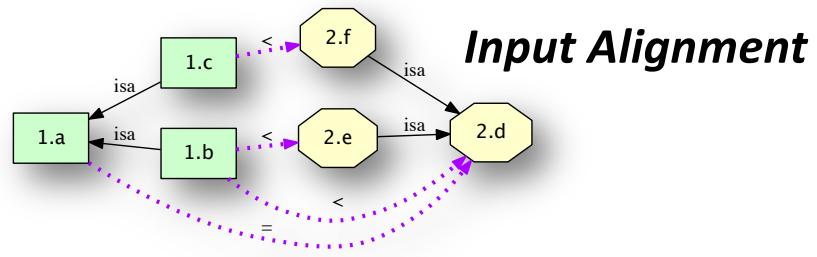
Taxonomy Alignment Problem (TAP)



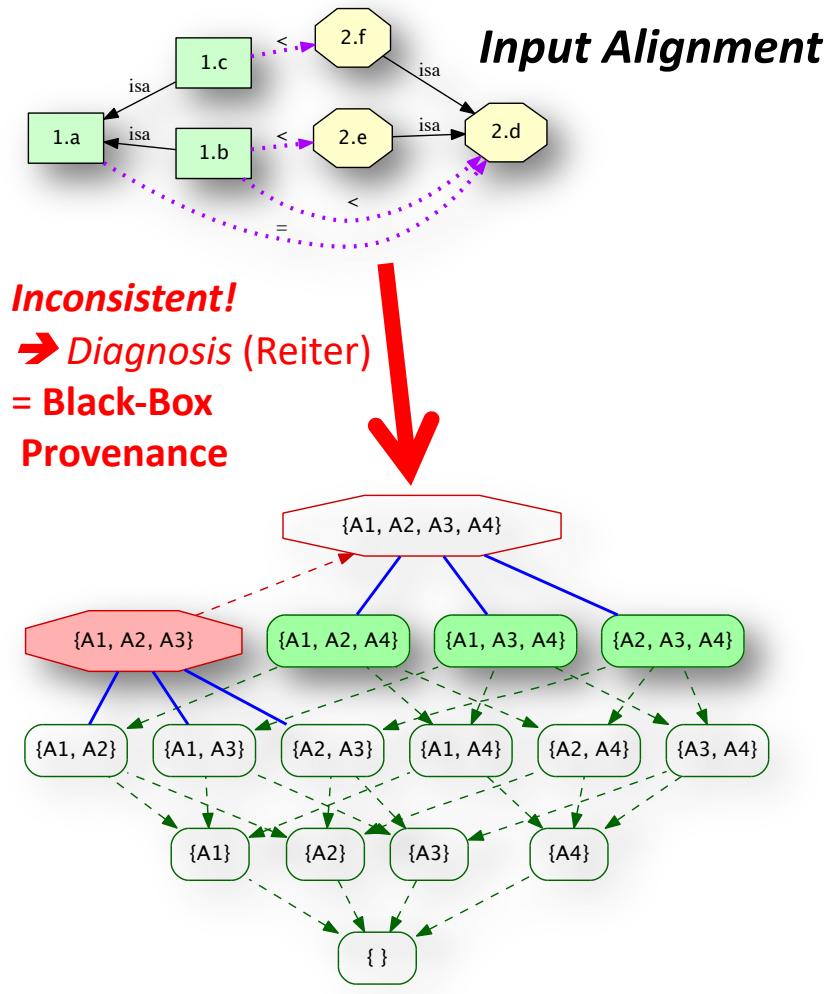
- **Given:**
 - Taxonomies T_1 , T_2
 - incl. constraints (coverage, disjointness)
 - Set of articulations (*alignment*) A
- **Find:**
 - **Combined (“merged”) taxonomy T_3** ($= T_1 + T_2 + A$)
 - Is it a taxonomy? Or a DAG?
 - Optional:
 - Final alignment (should be minimal)



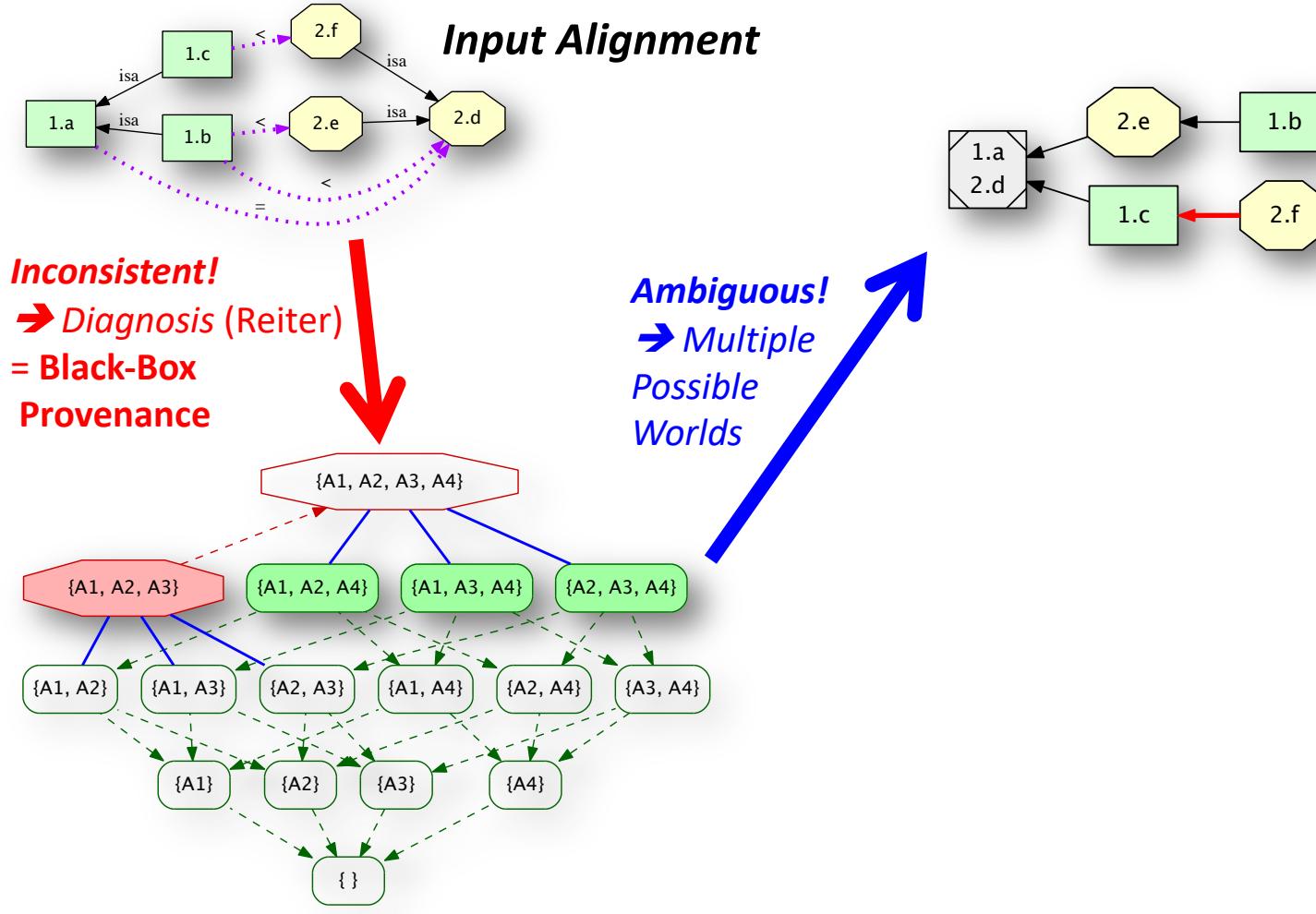
TAP: Possible Outcomes



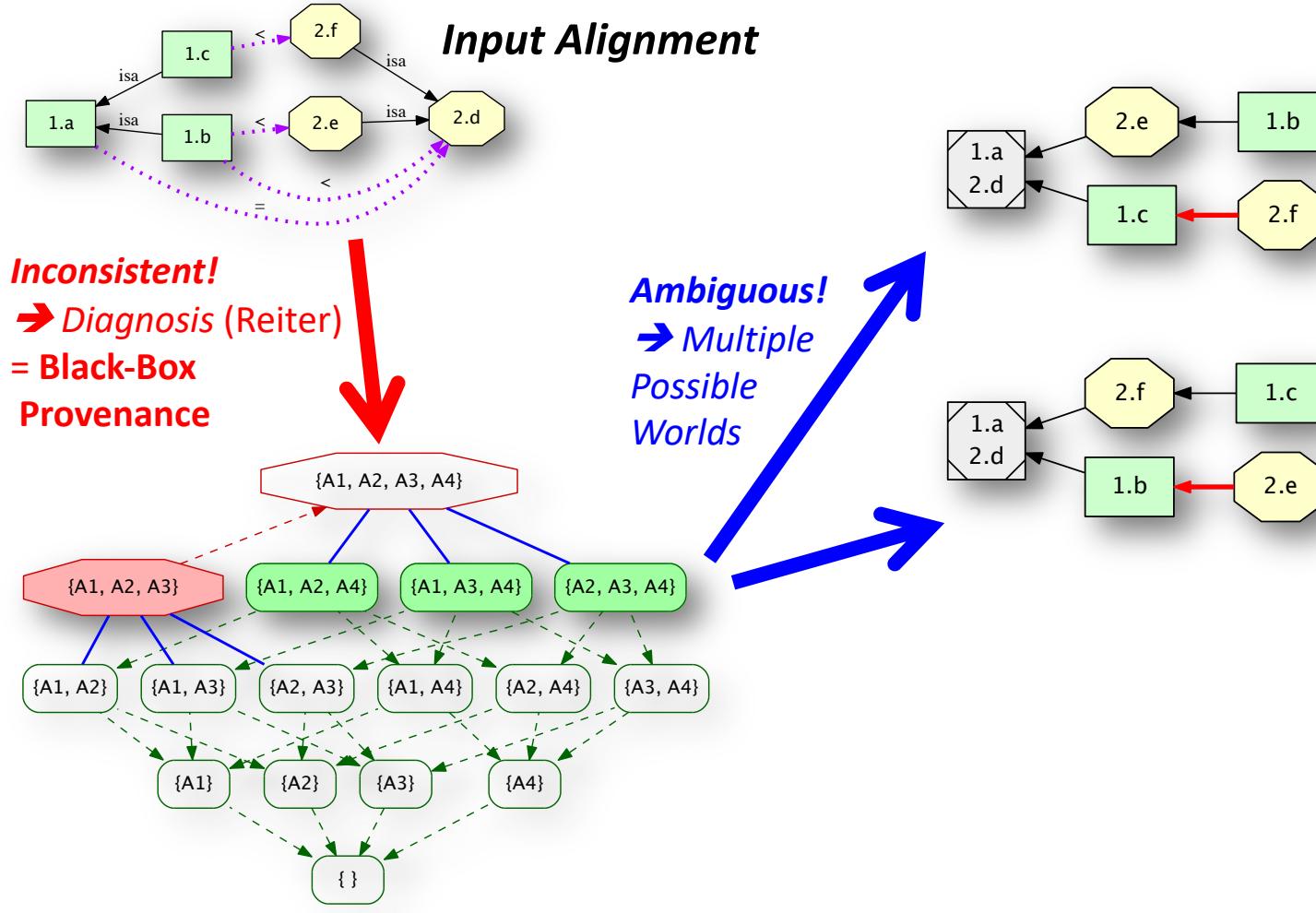
TAP: Possible Outcomes



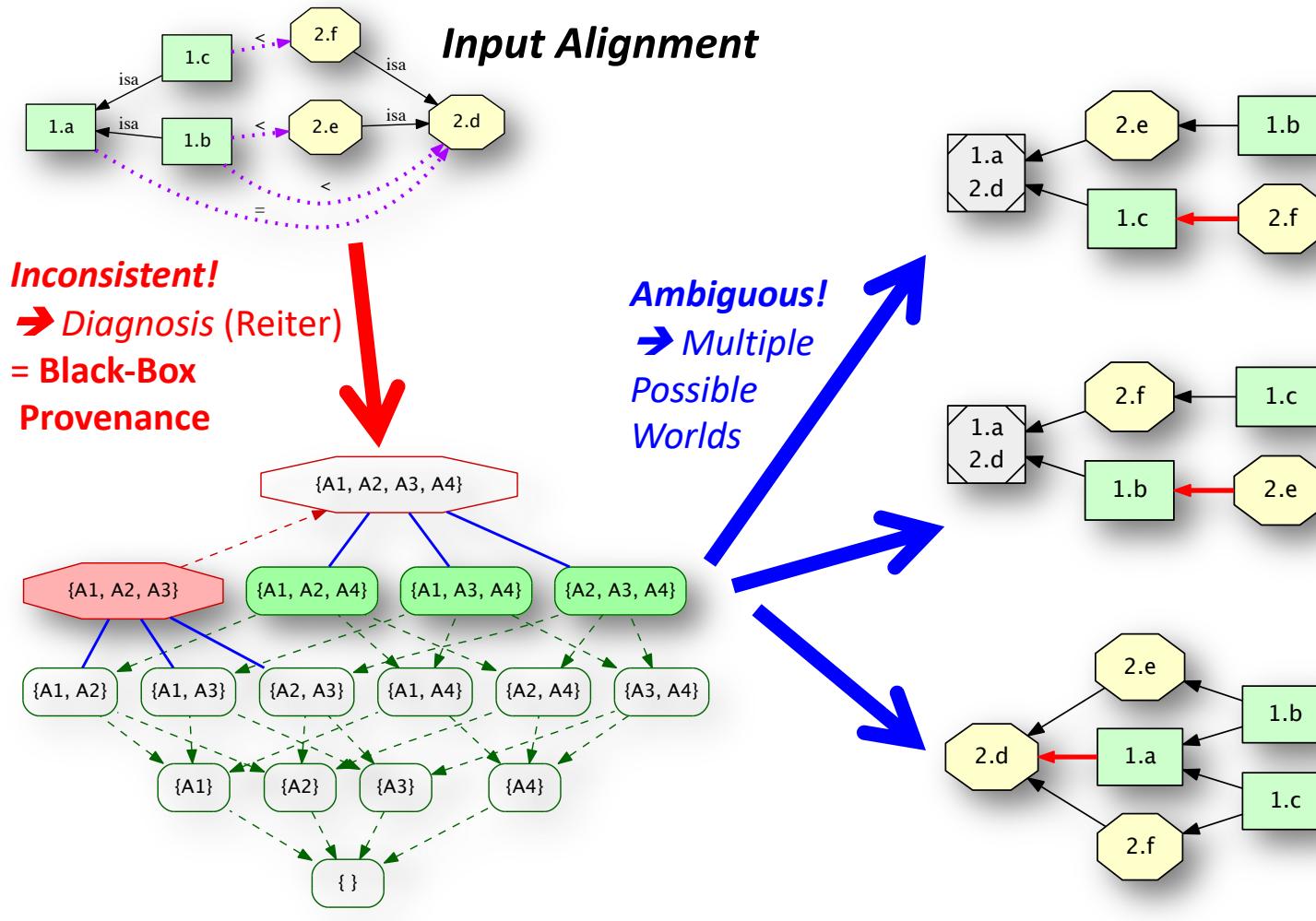
TAP: Possible Outcomes



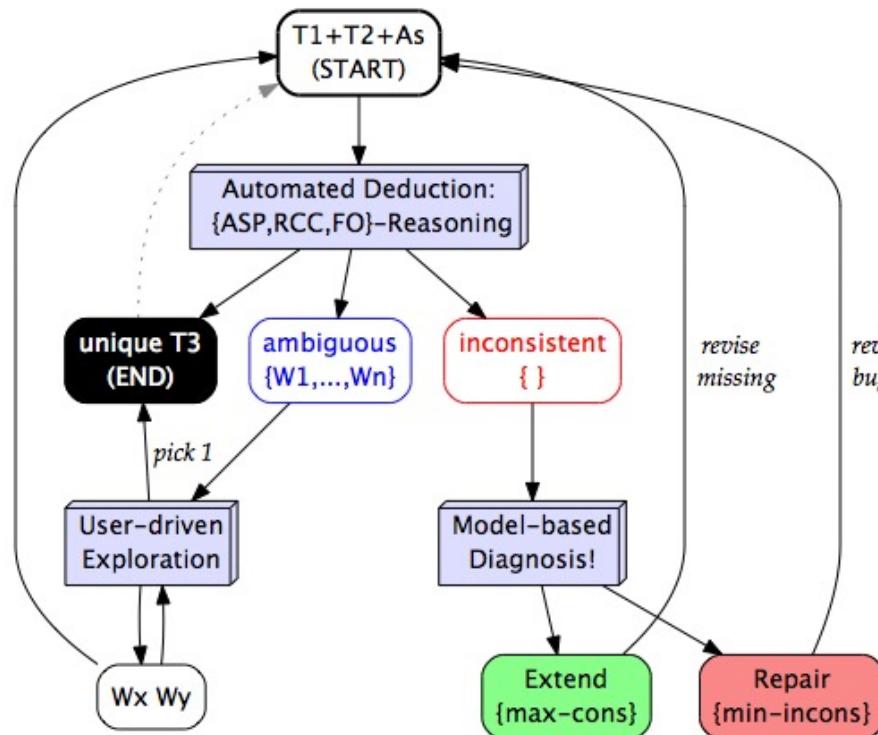
TAP: Possible Outcomes



TAP: Possible Outcomes



Euler/X Toolkit and Workflow



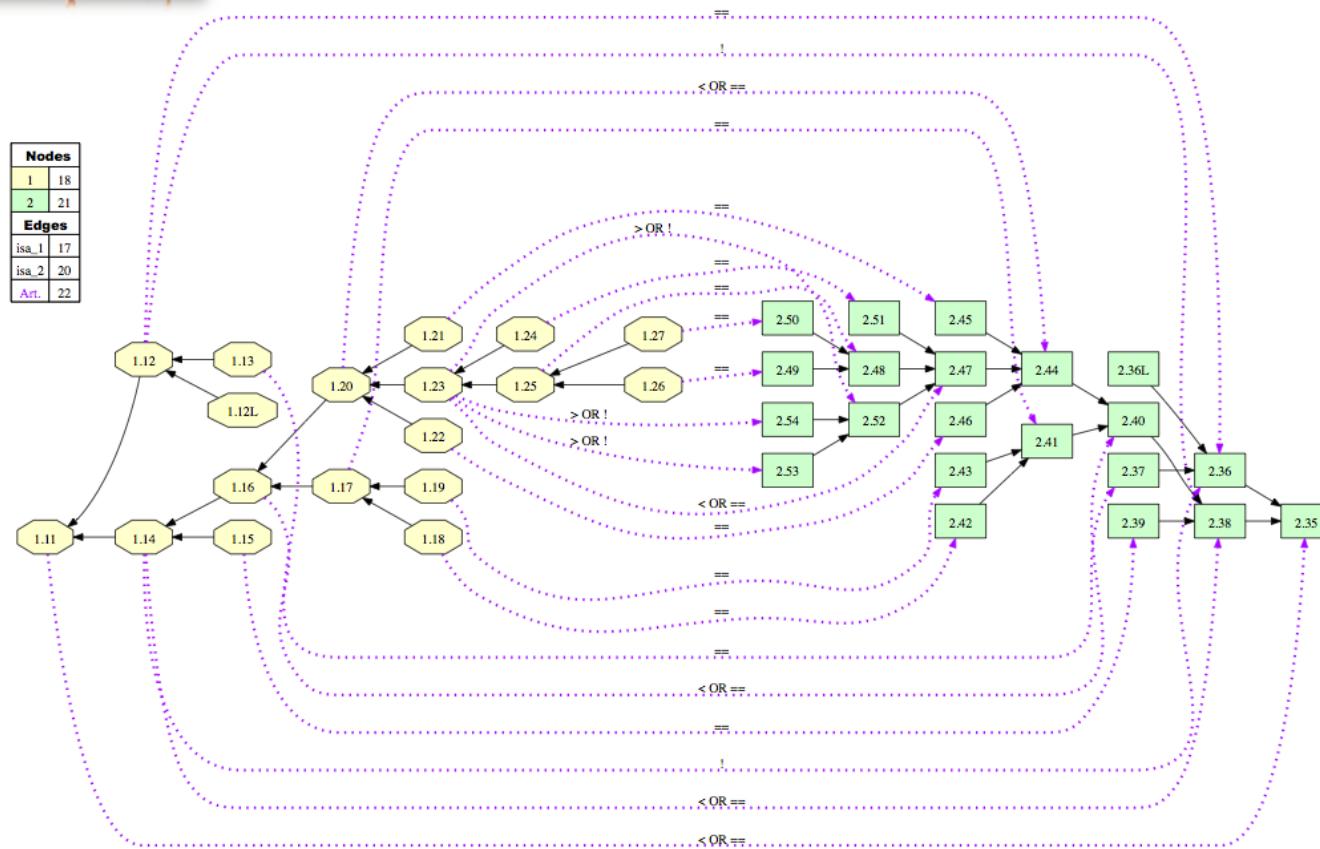
- FO reasoning about taxonomies (MFOL)
- Earlier: **CleanTax**
 - Prover9/Mace4
- Now: **Euler**
 - ASP Reasoners (DLV, Clingo)
 - Specialized reasoners (PyRCC)
 - ...
 - X = ASP, RCC, ...



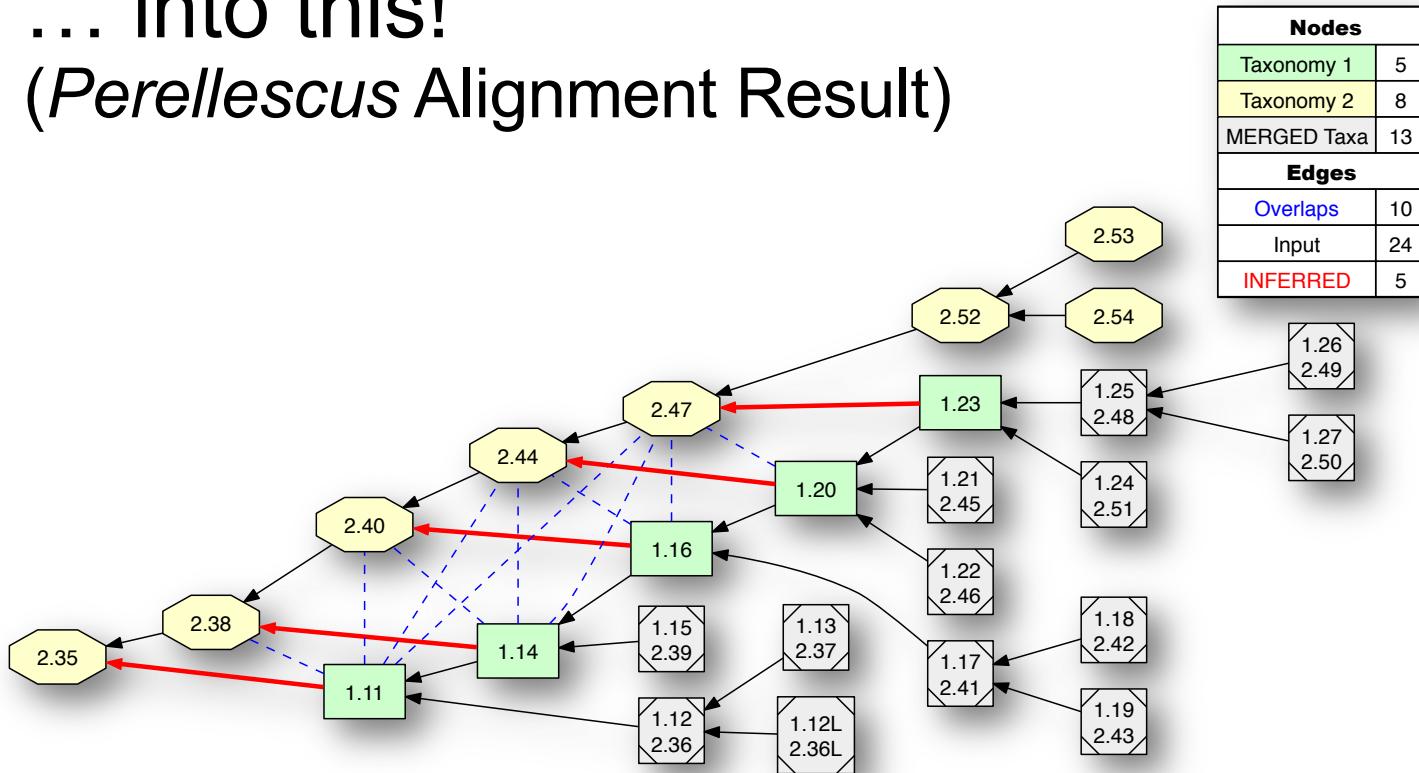
8

Real-world examples: Turn this

1



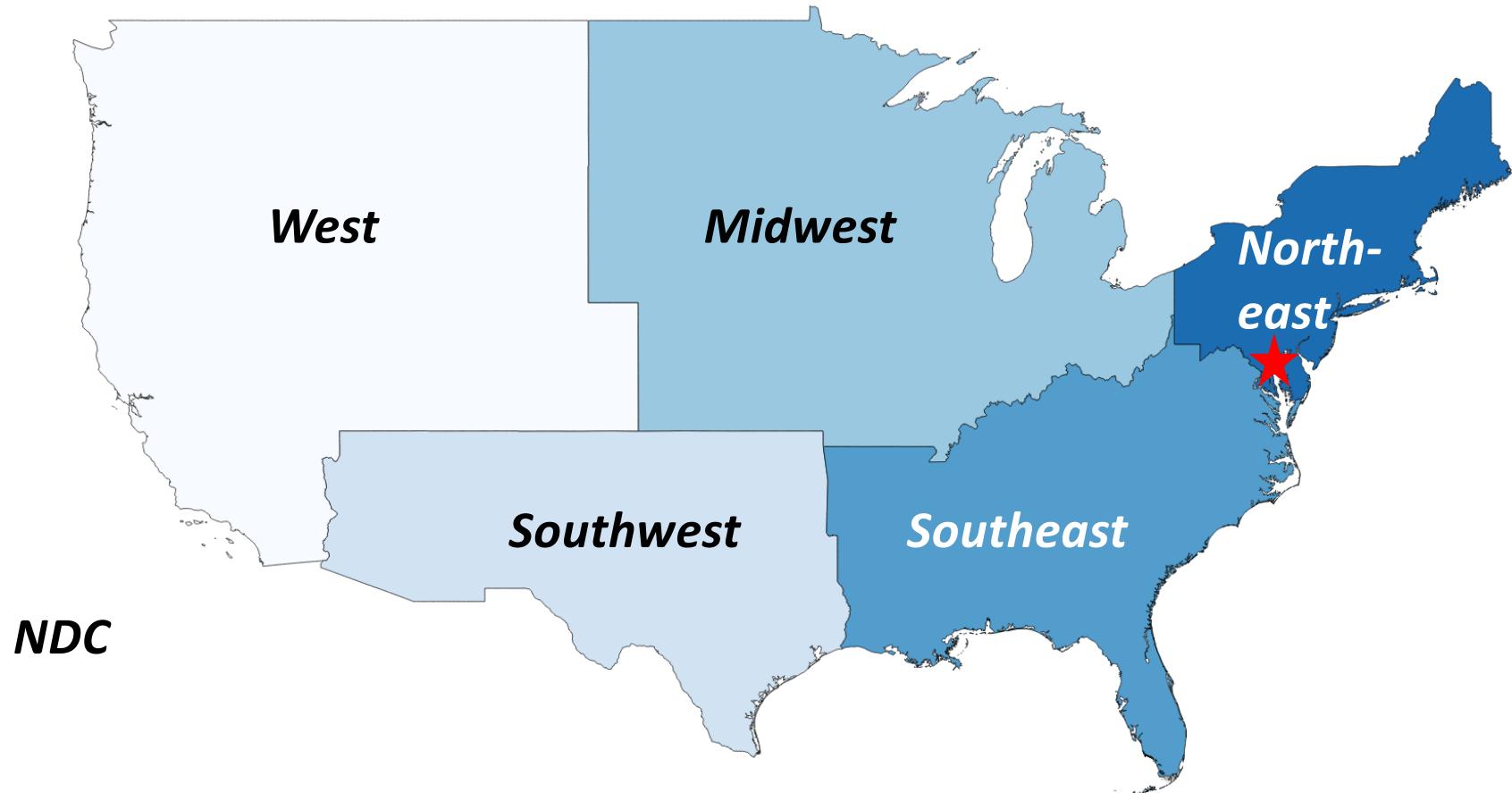
... into this! (*Perellescus* Alignment Result)

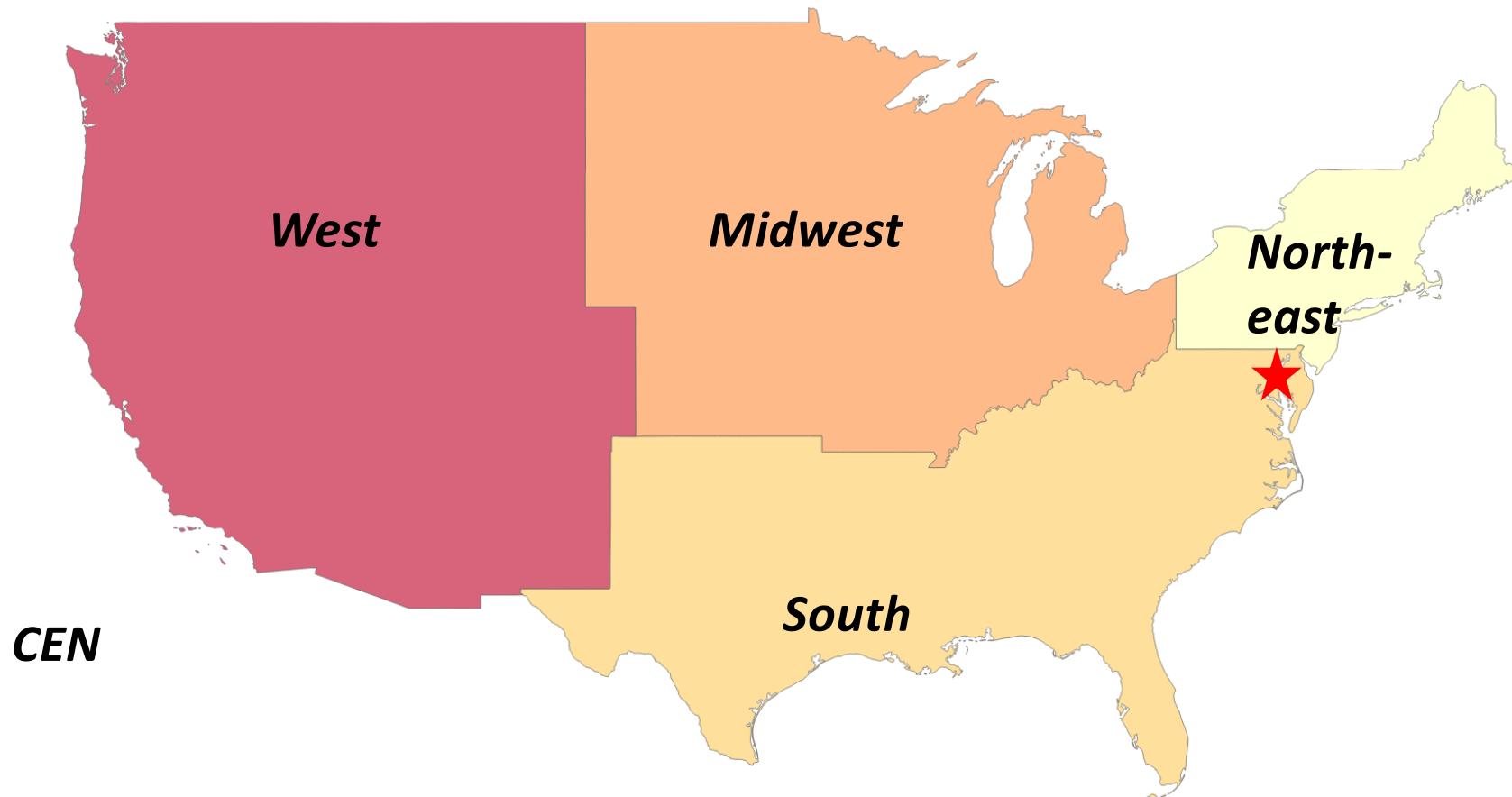


- $T_3 := T_1$ and T_2 are “merged”
 - Blue dashed: overlaps → resolve via “zoom-in view”

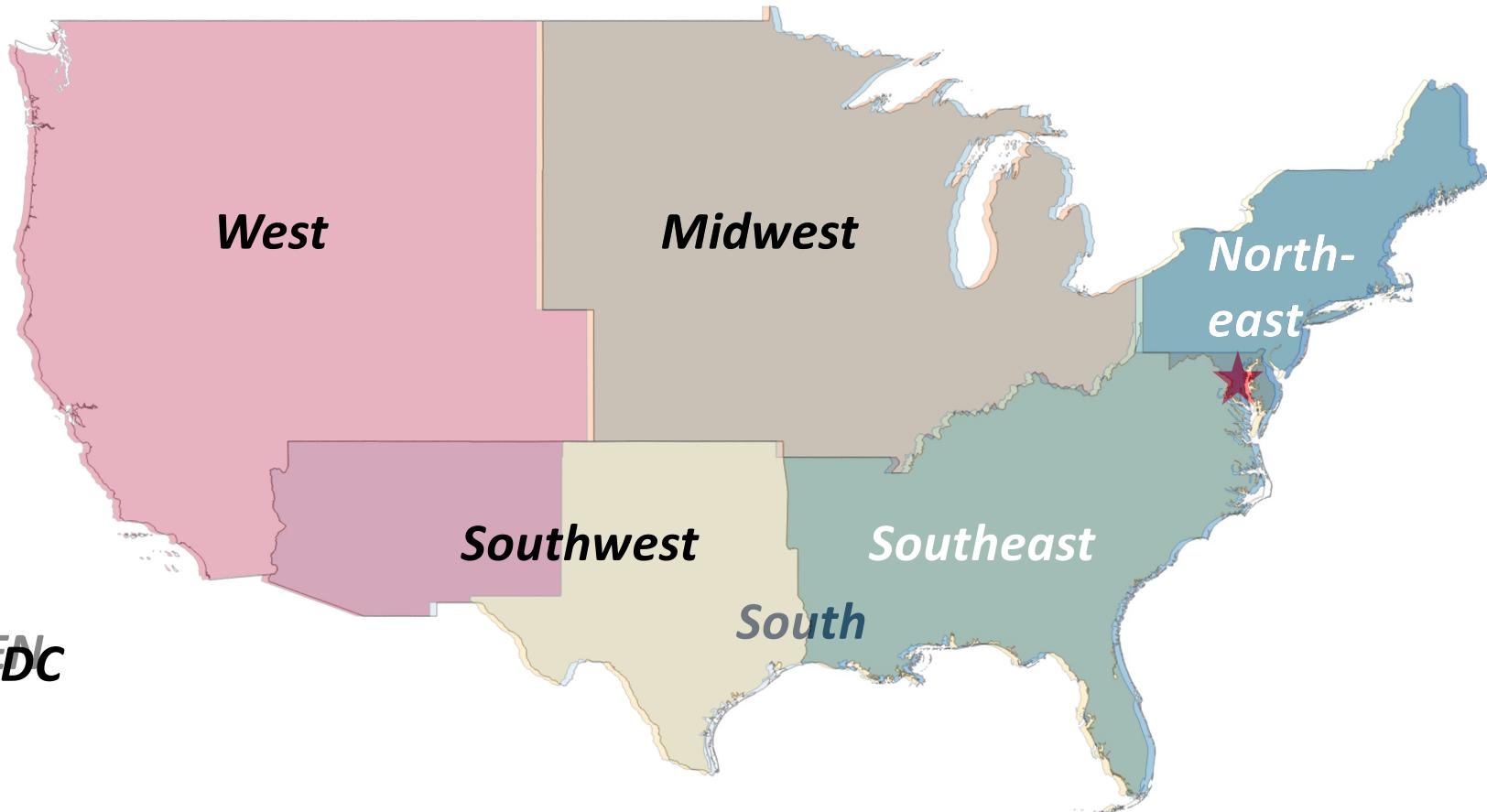
CEN-NDC-Regions

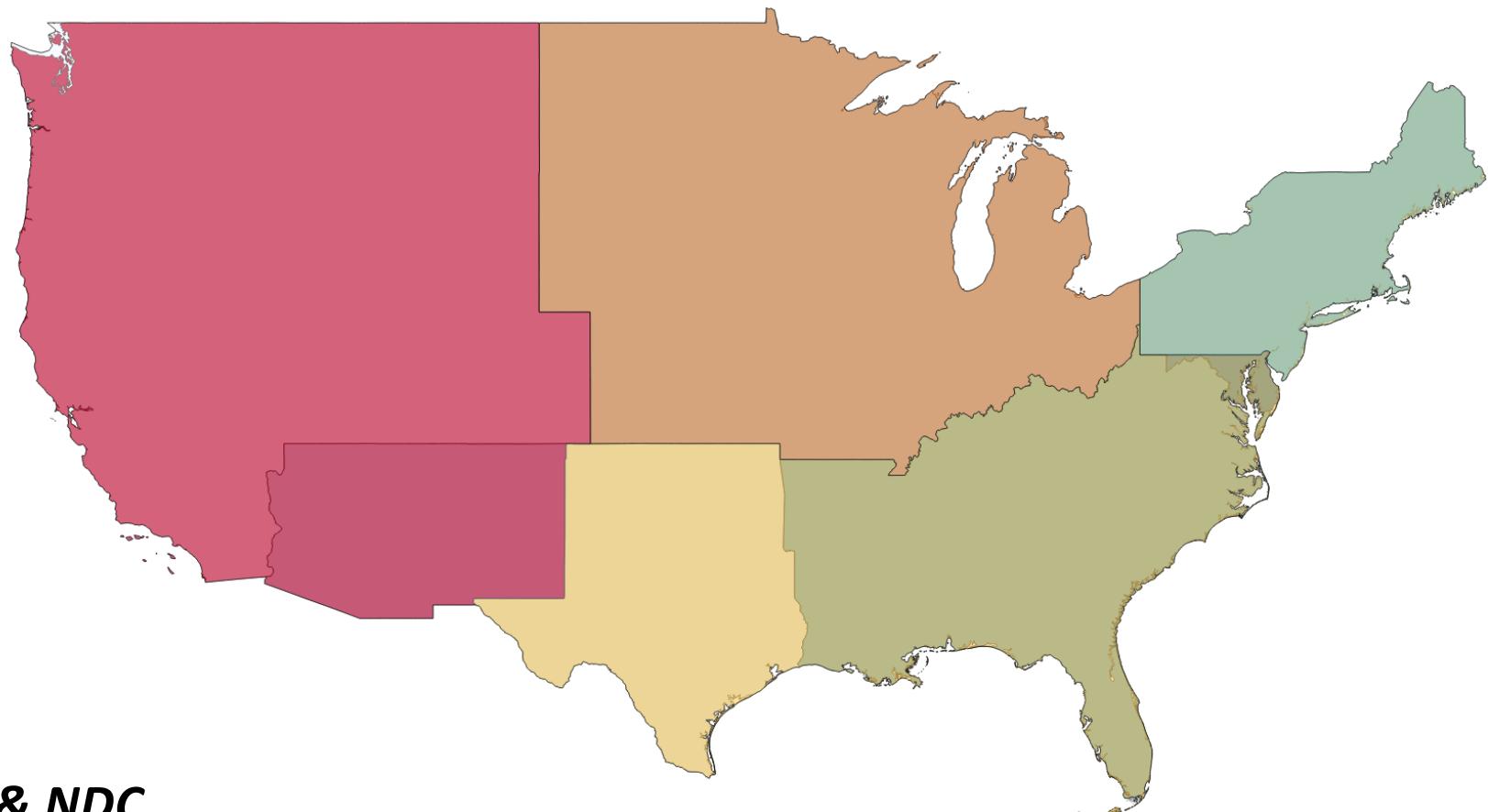
Alignment



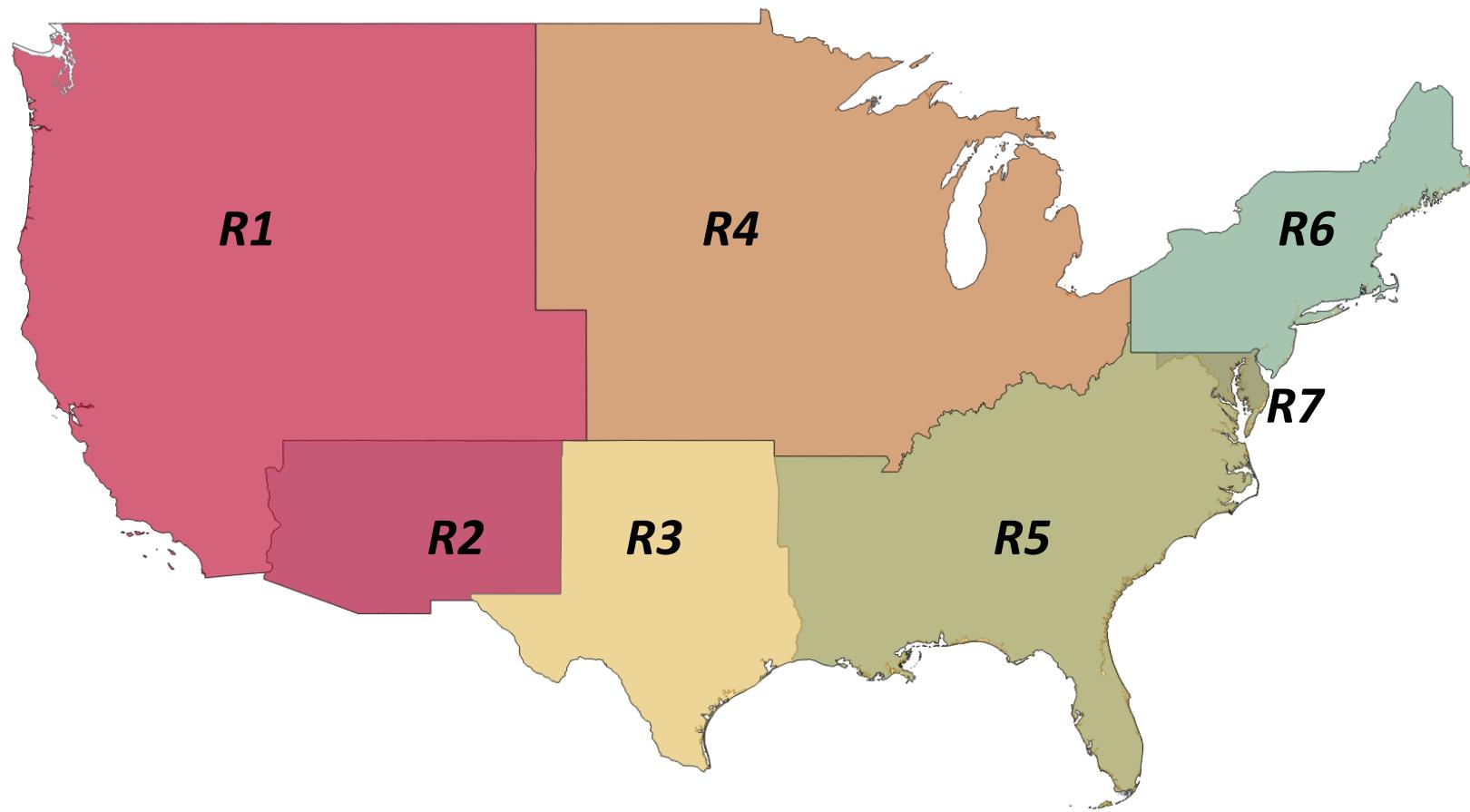


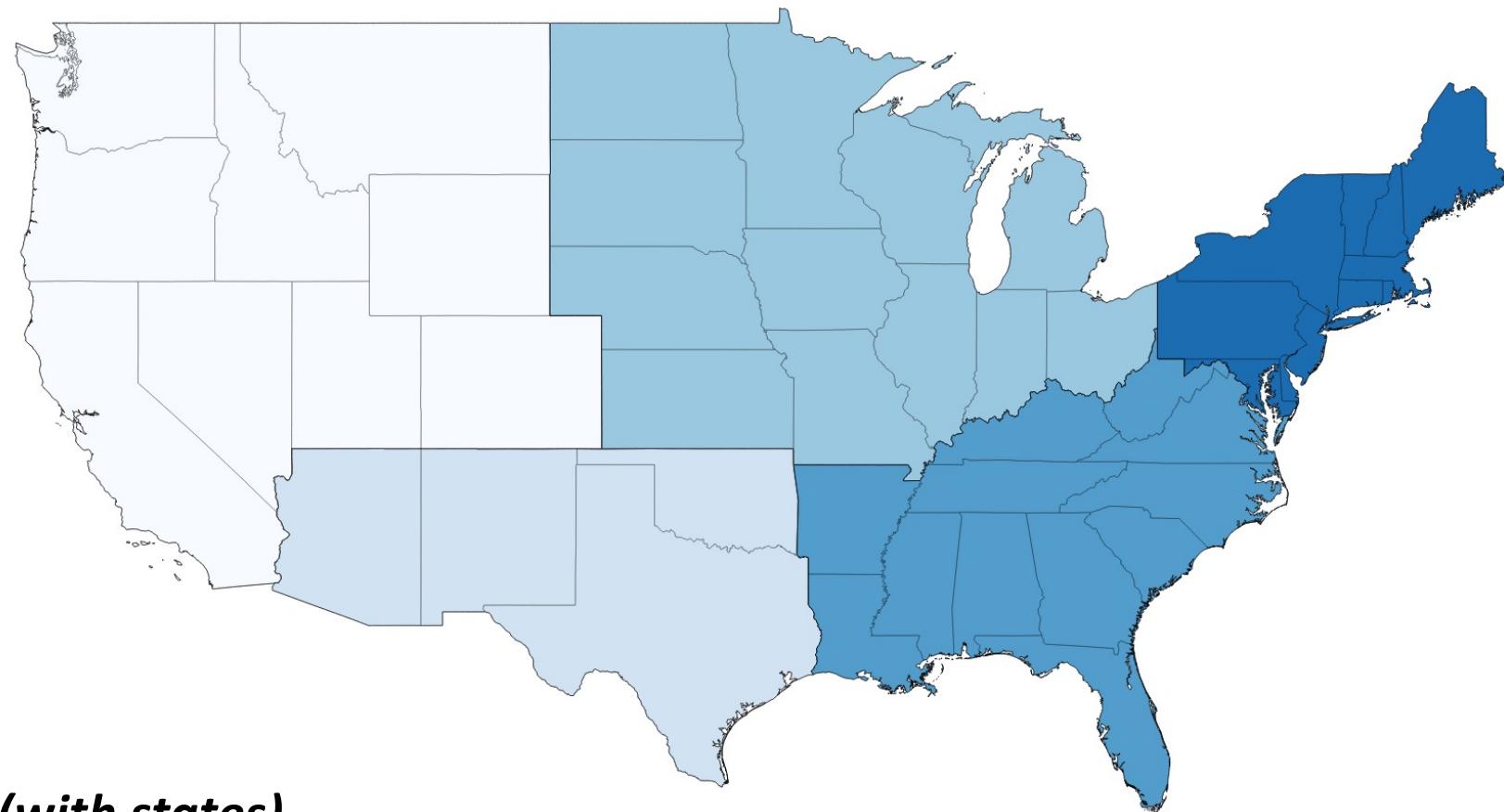
**G
EN
NDC**



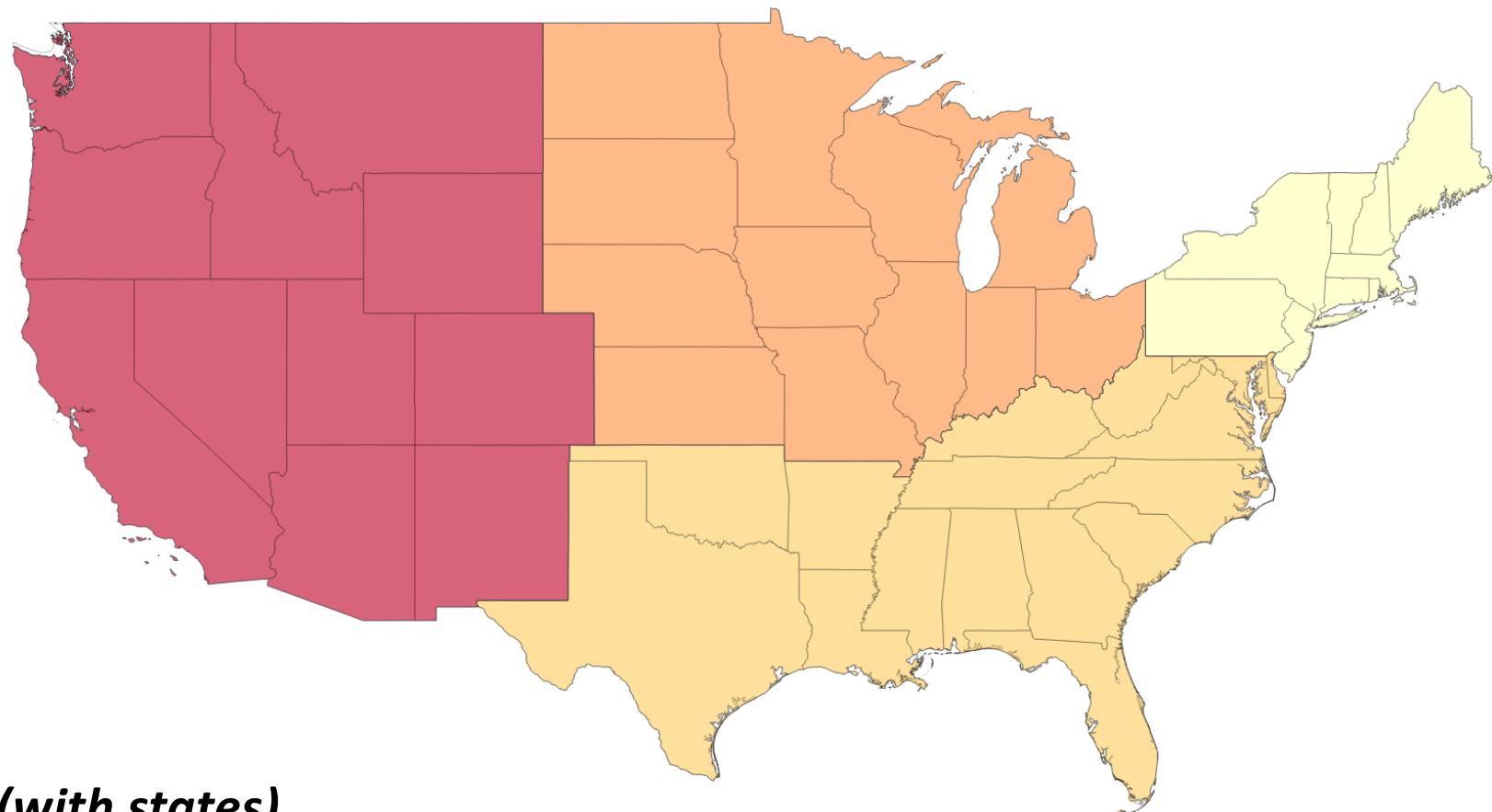


CEN & NDC

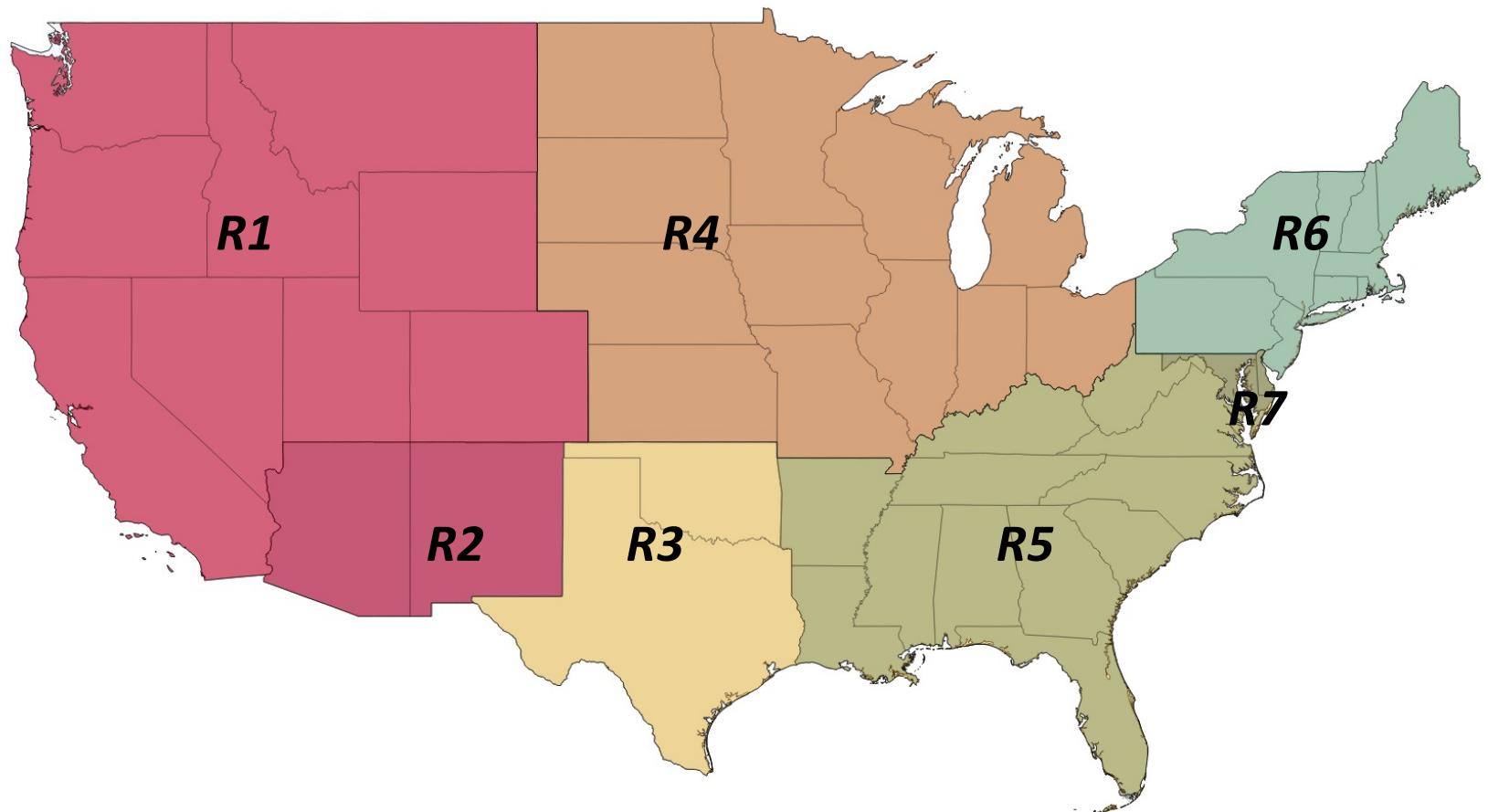


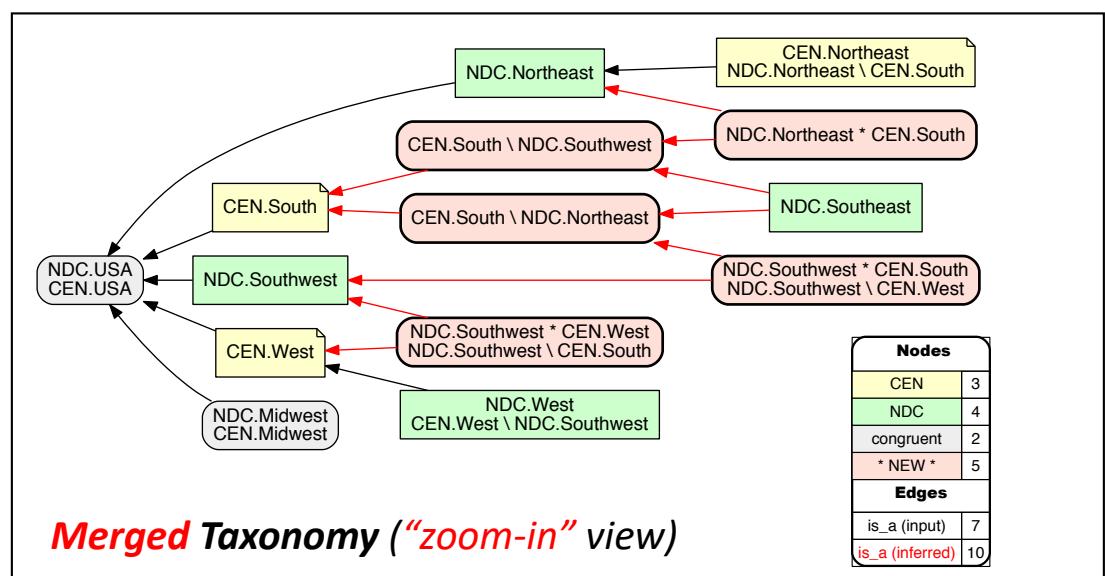
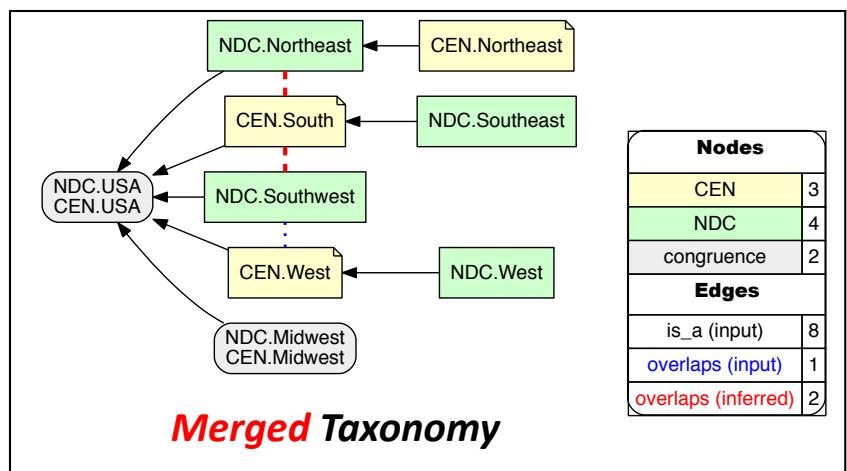
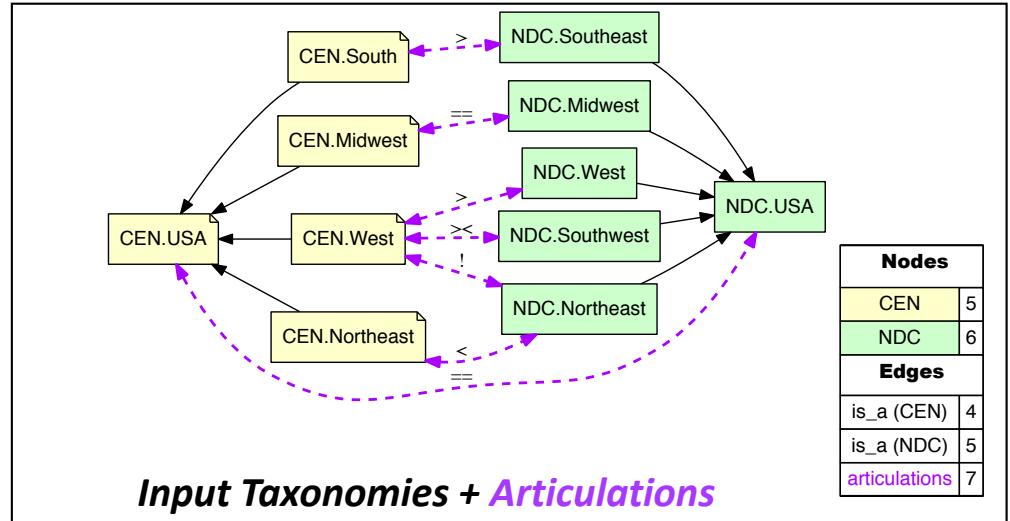
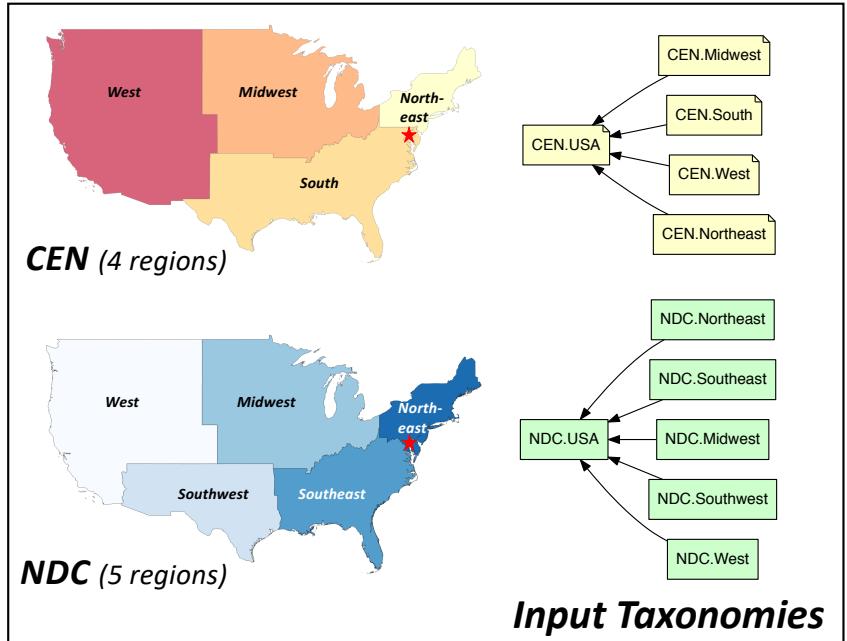


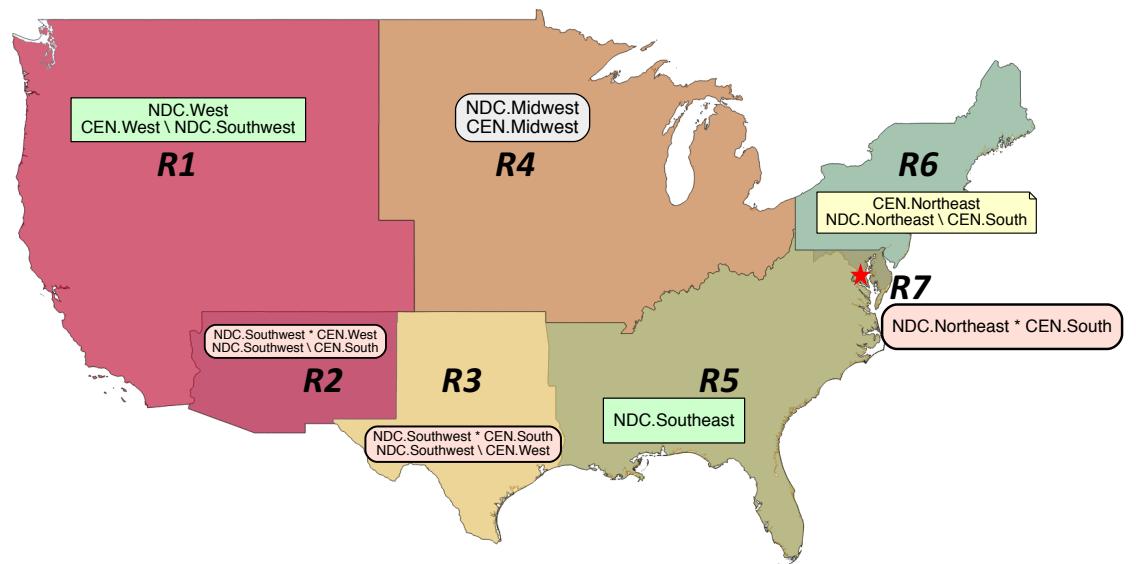
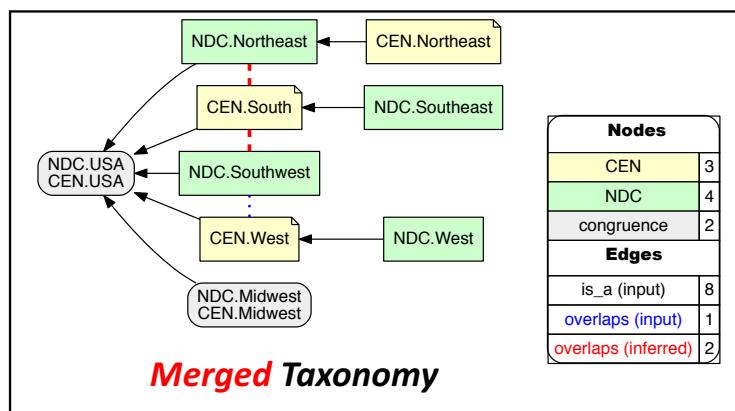
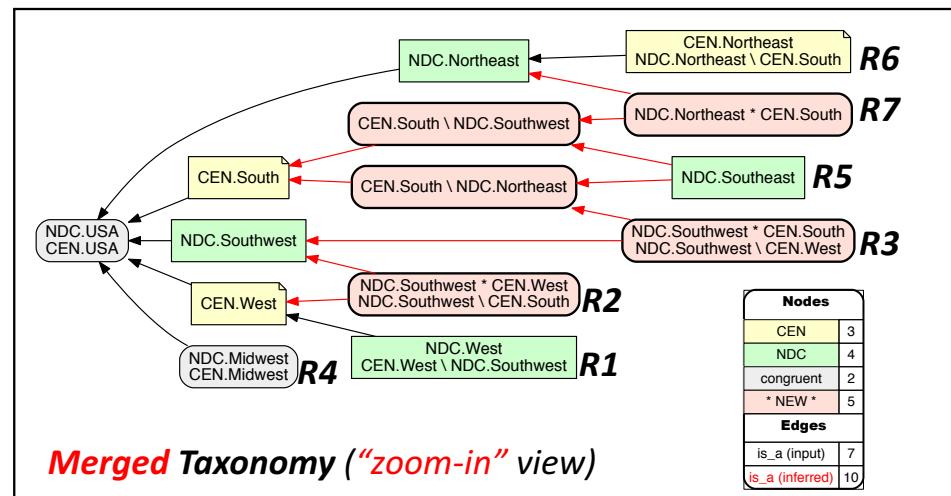
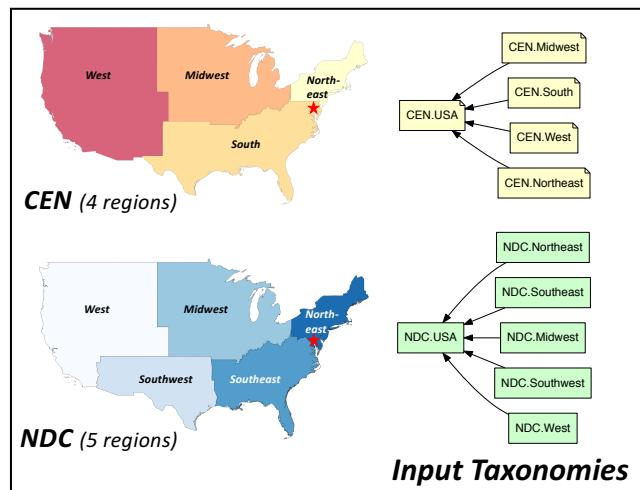
NDC (with states)

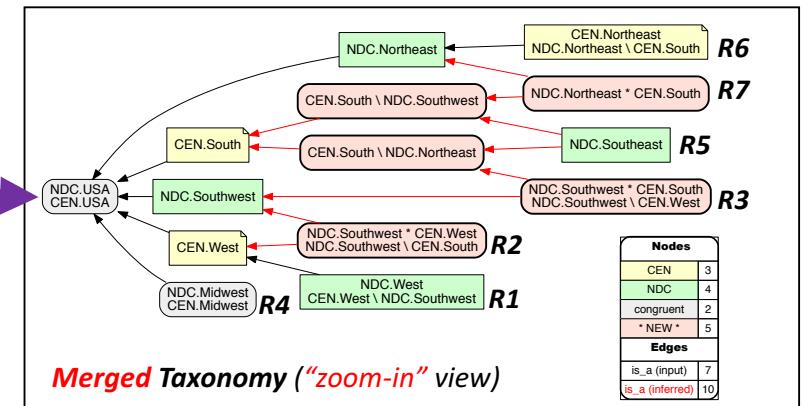
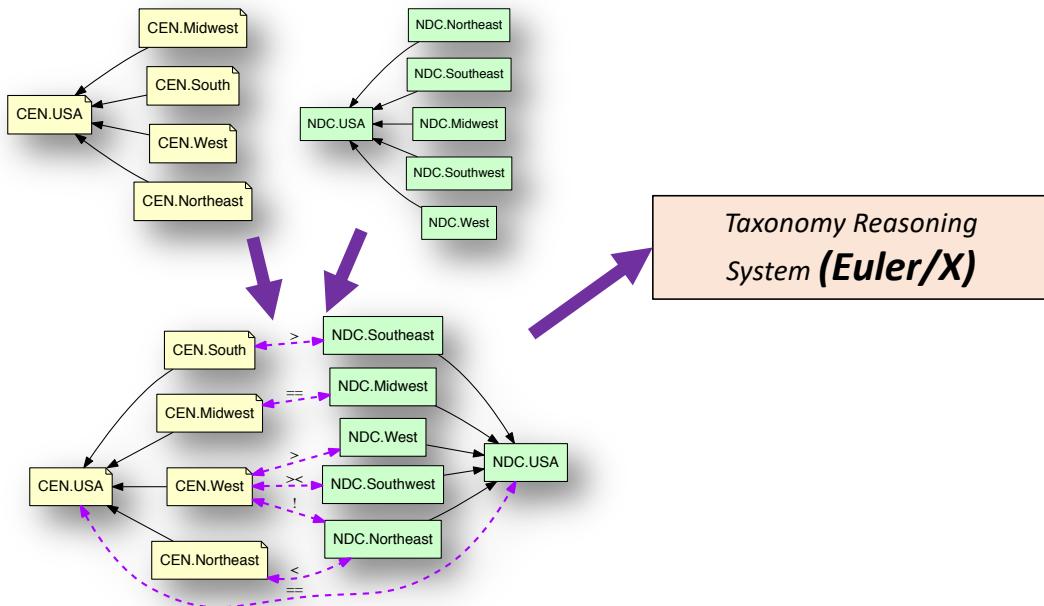
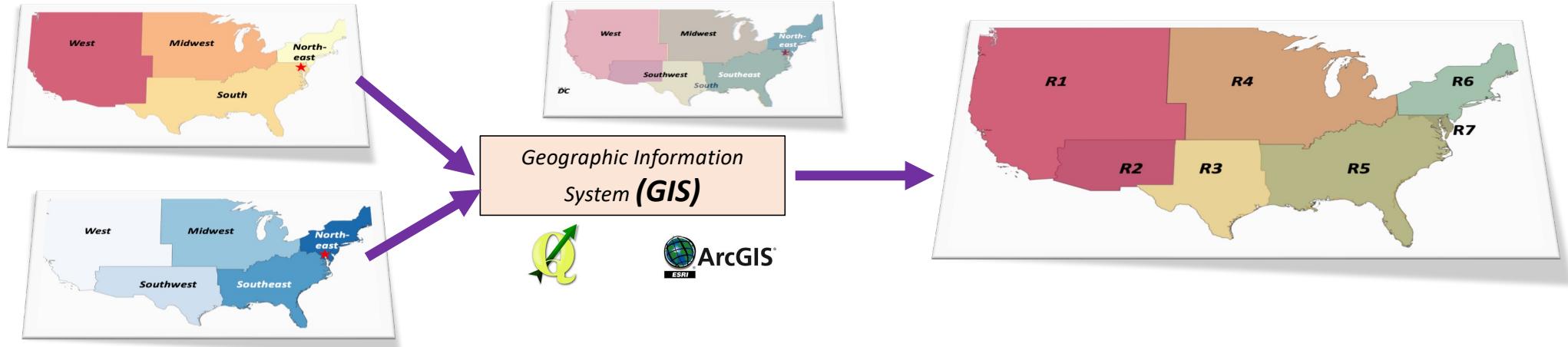


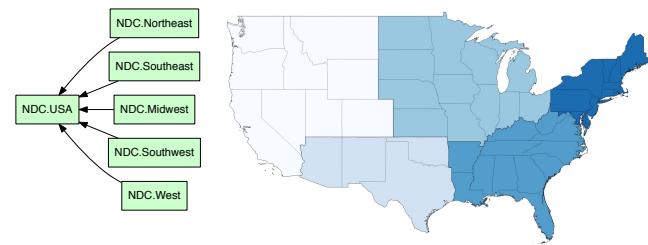
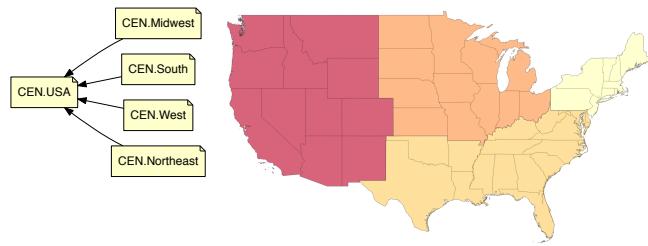
CEN (with states)



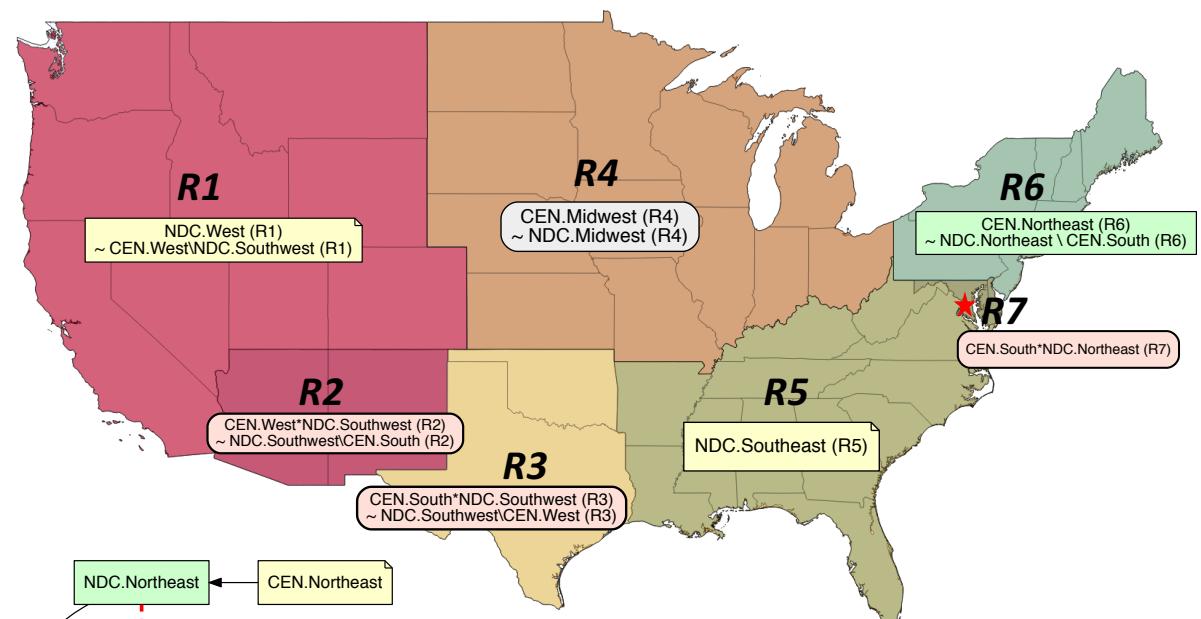
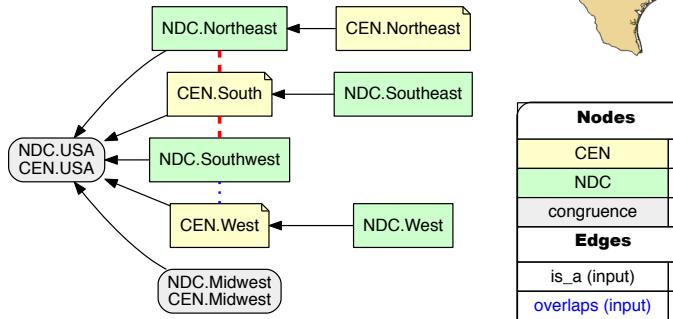
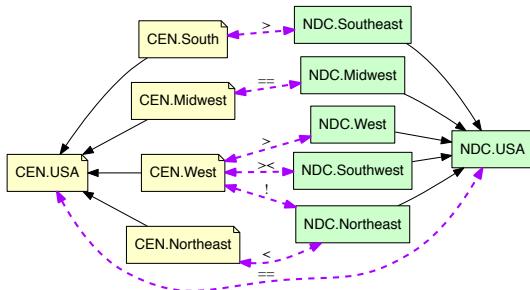






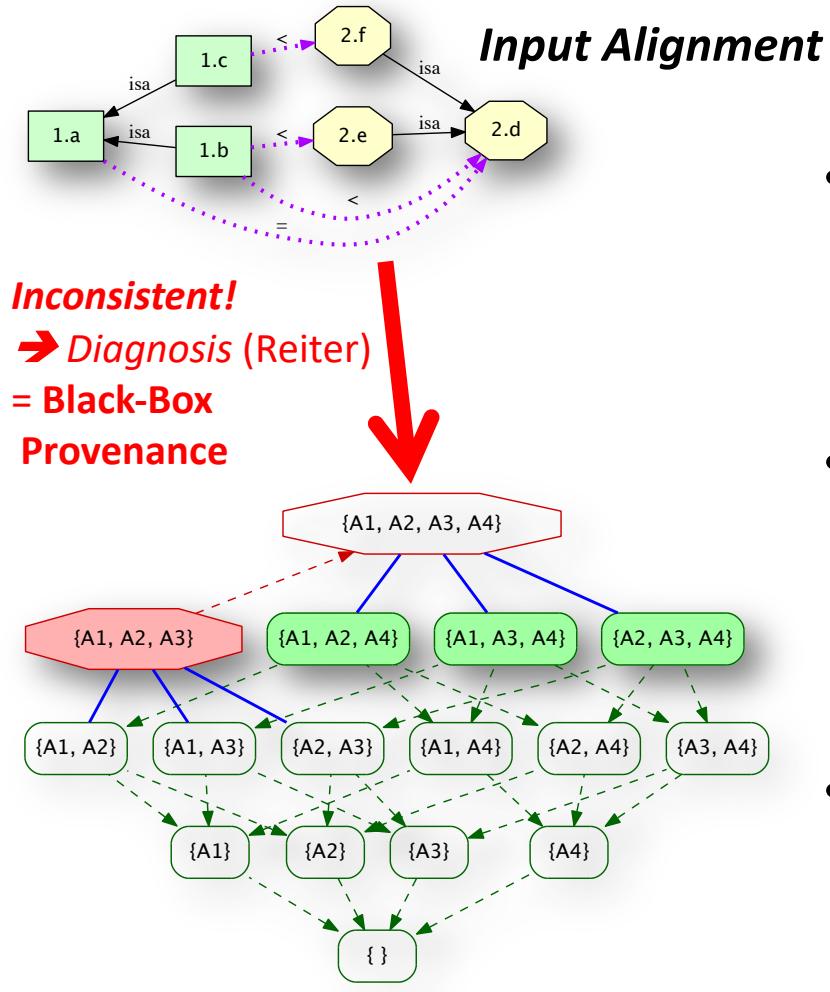


Nodes	
CEN	5
NDC	6
Edges	
is_a (CEN)	4
is_a (NDC)	5
articulations	7



Nodes	
CEN	3
NDC	4
Edges	
is_a (input)	8
overlaps (input)	1
overlaps (inferred)	2

Possible Outcome: Inconsistency!

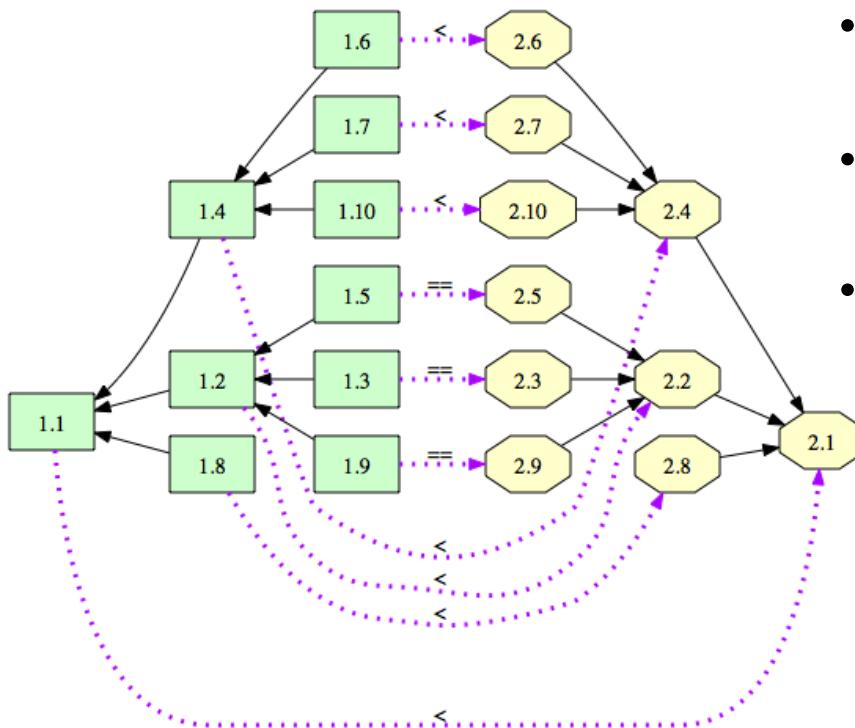


- Need to **debug** the input articulations → (black-box) diagnosis!
- Focus:
 - How do we **efficiently** compute the diagnostic lattice?
- Also:
 - How to **visualize..**

Nodes	
1	10
2	10
Edges	
isa_1	9
isa_2	9
Art.	10

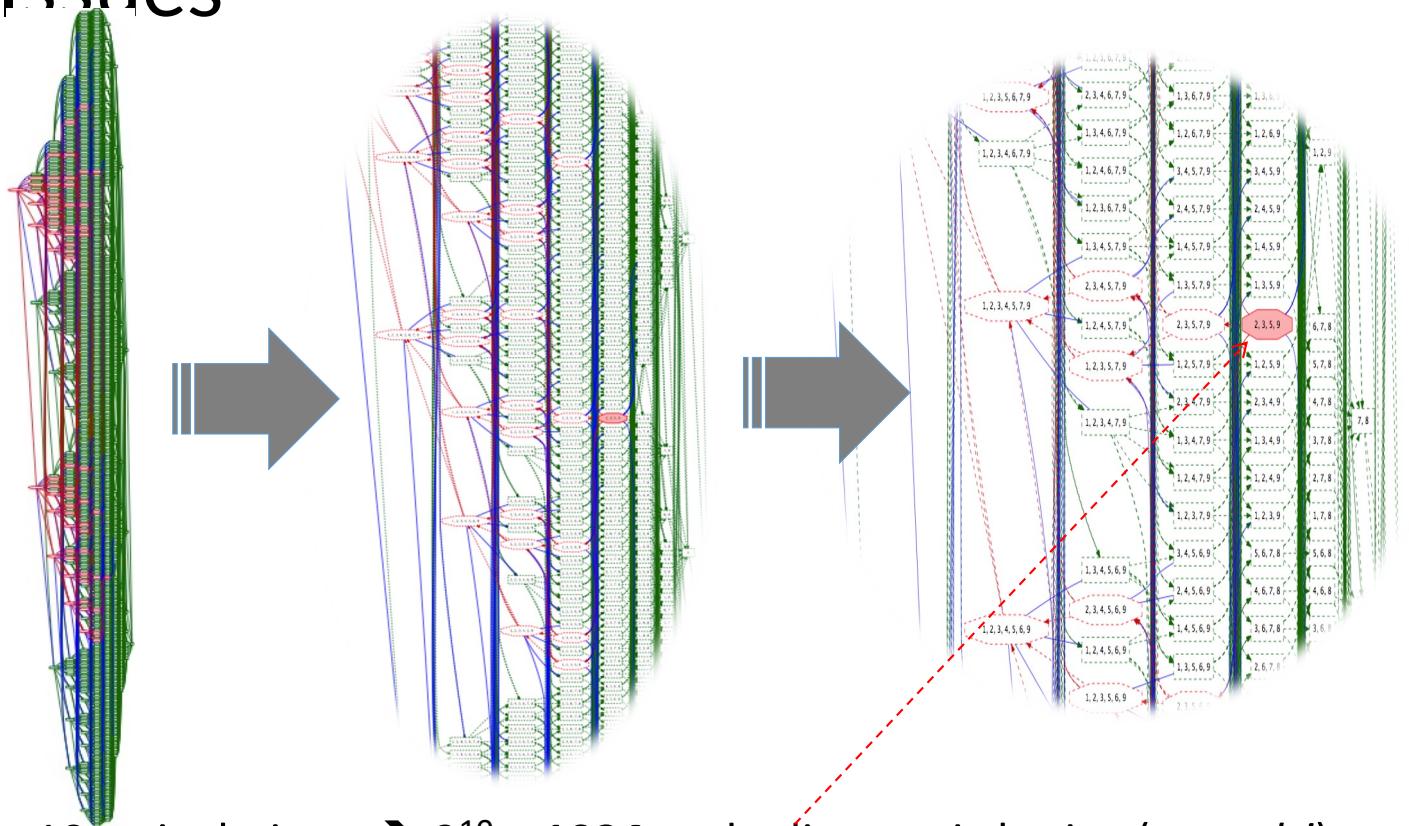
Example Instance

(from synthetic benchmark suite)



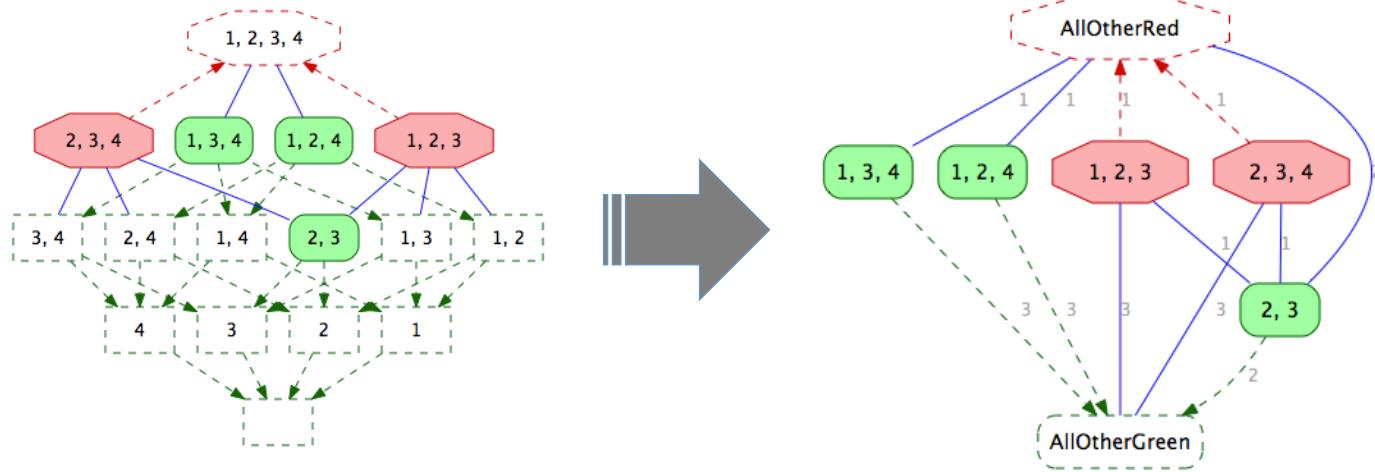
- Here: $N = 10$ taxa in T_1 , T_2
- Euler/X finds:
inconsistent!
- → diagnostic lattice of $2^{10} = 1024$ nodes
 - Find **minimal inconsistent subset (MIS)**
 - maximal consistent subset (**MCS**) ..
 - show to user!

Visualizing Diagnoses: Scalability Issues



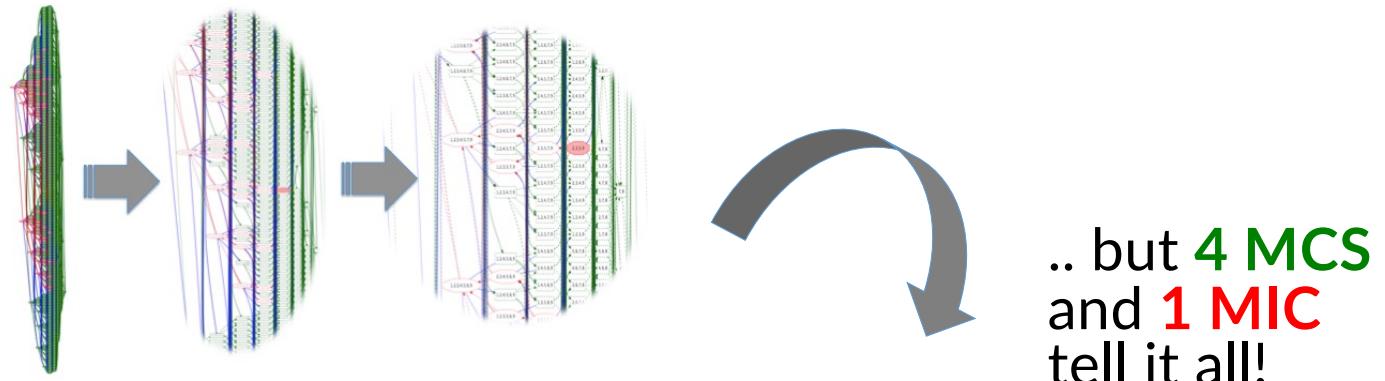
N = 10 articulations $\rightarrow 2^{10} = \mathbf{1024}$ node diagnostic lattice (... ouch!)
... but only one MIS ...

Better Idea: Just show MIS, MCS

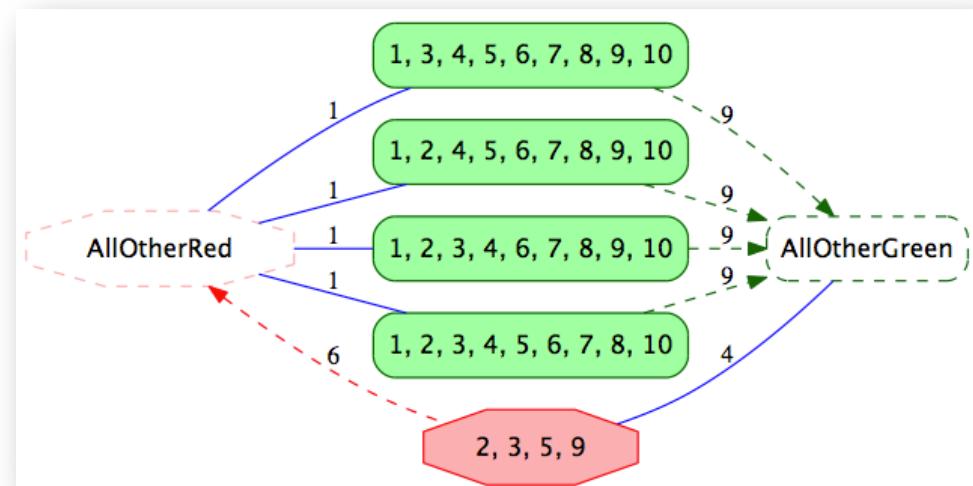


$N = 4$ articulations $\rightarrow 2^4 = 16$ node diagnostic lattice,
but 3 **MCS** and 2 **MIS** are enough!

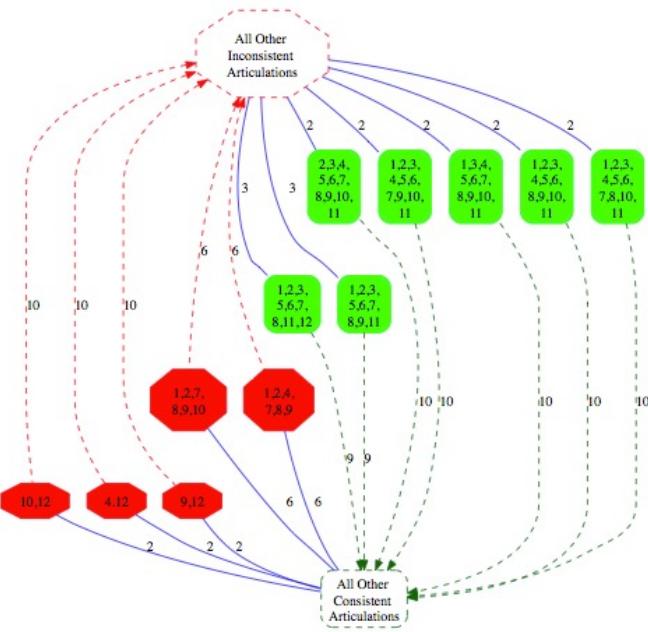
Visualizing Diagnoses: Focusing on MIS (and MCS)



1024 node lattice



Visualizing Diagnoses

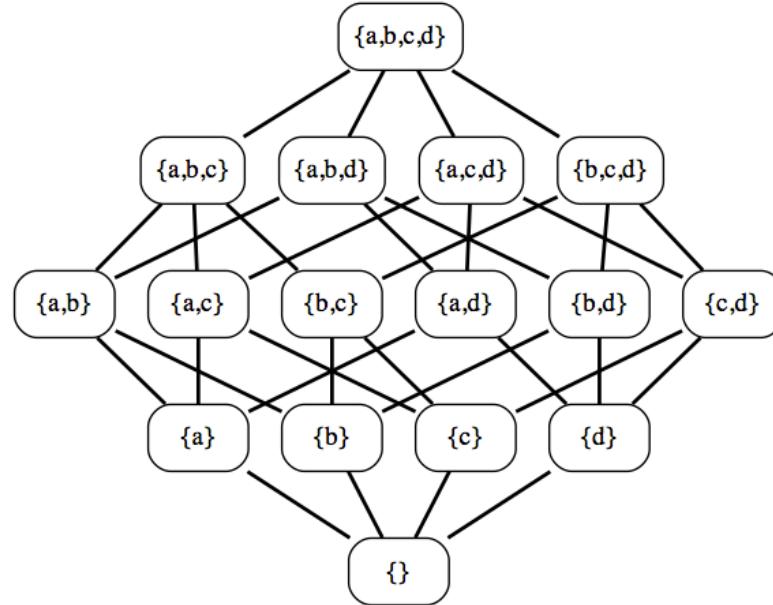


Example from
paper: N=12 →
4096 nodes

.. but **7 MCS** and
5 MIC tell it all!

Fig. 7. \mathcal{MIS} (Octagon) and \mathcal{MCS} (Rounded Box) for Example 2. All other non-minimal inconsistent subsets and non-maximal consistent subsets are collapsed as “clouds”, the labels of edges show the path length from $\mathcal{MIS}/\mathcal{MCS}$ to the top/bottom of the lattice.

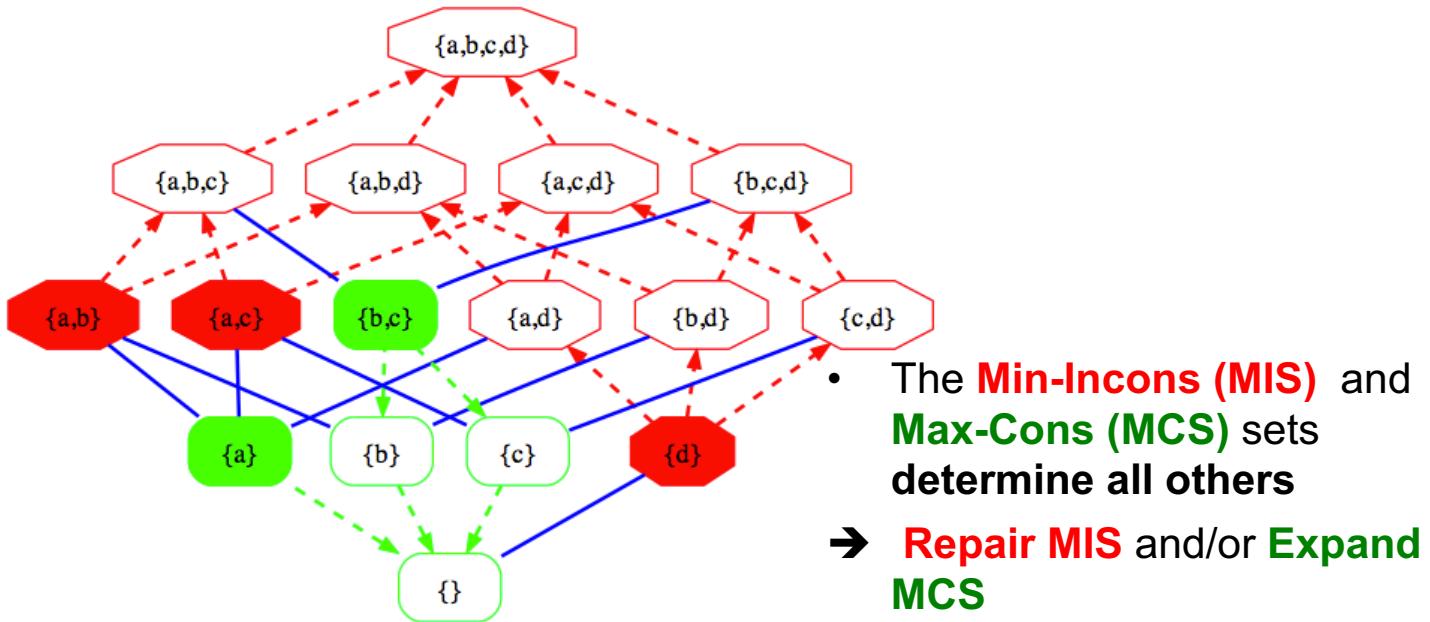
Black-Box Inconsistency Analysis (Diagnostic Lattice)



What happens if you can't have all (here: 4) articulations together?

- Then:
 - Repair: find & revise **minimal inconsistent** subsets (**Min-Incons**)
 - Expand: find **maximal consistent** subsets (**Max-Cons**) & revise *outs*

Inconsistency Analysis (Diagnostic Lattice)



- **Black-box Analysis** (Hitting Set algo.) yields a **Diagnosis** (lattice)
 - for $n=4$ articulations, there are 168 possible diagnoses
 - depending on expected “red/green areas” → explore space differently
- $|\text{articulations}| = n \rightarrow$

$$|\text{possible diagnoses}| = |\text{monotonic Boolean functions}|$$

$$= \text{Dedekind Number } (n): 2, 3, 6, 20, \mathbf{168}, 7581, 7828354, \dots$$

Improving Diagnosis

- Reiter’s “black-box” (model-based) diagnosis helps debug the articulations
 - Limited scalability (inherent complexity)
 - But every bit helps:
 - Hitting Set Algorithm (“logarithmic extraction”)
 - Our idea:
 - Exploit “white-box” reasoning information
- ➔ **RULES** to the rescue

Key Idea: exploit white-box info

- We use Answer Set Programming (ASP) to solve Taxonomy Alignment Problem (TAP)
- Inconsistency = “False” is derived in the head:
False :- <denial of integrity constraint>
- Apply **provenance trick** from databases ☺
 - What articulations contribute to a derivation of “False” ?
 - **Eliminate those that don’t!**
- → an example of reusing inferences **across** separate black-box tests!

The Provenance “Trick”

$$H_1(\bar{Y}, r_1 \otimes (P_1 \otimes \dots \otimes P_n)) :- B_1(\bar{X}_1, P_1), B_2(\bar{X}_2, P_2), \dots, B_n(\bar{X}_n, P_n).$$

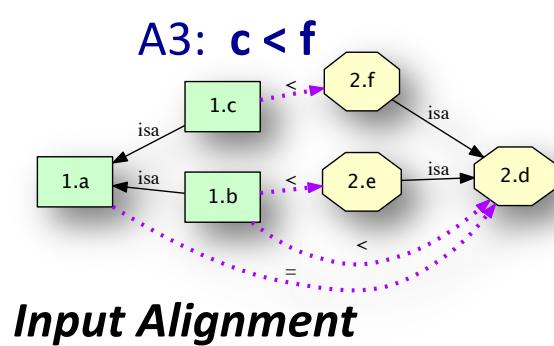
For any constraint rule without head predicate (i.e., `false` is the head):

$$r_2 : \text{false} :- B_1(\bar{X}_1), B_2(\bar{X}_2), \dots, B_n(\bar{X}_n).$$

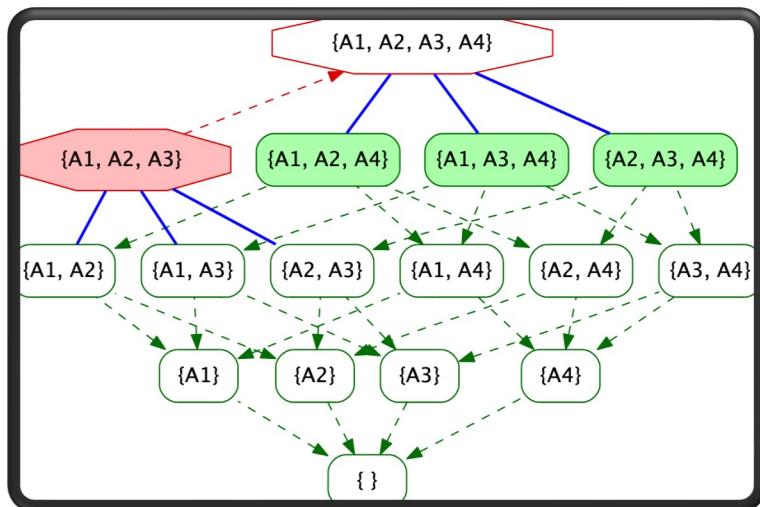
We rewrite it to a constraint with head predicate `NOK` where P_i is the provenance of $B_i(\bar{X}_i)$ for $1 \leq i \leq n$ and `NOK` stands for “Not OK”, i.e. inconsistency:

$$\text{NOK}(r_2 \otimes (P_1 \otimes \dots \otimes P_n)) :- B_1(\bar{X}_1, P_1), B_2(\bar{X}_2, P_2), \dots, B_n(\bar{X}_n, P_n).$$

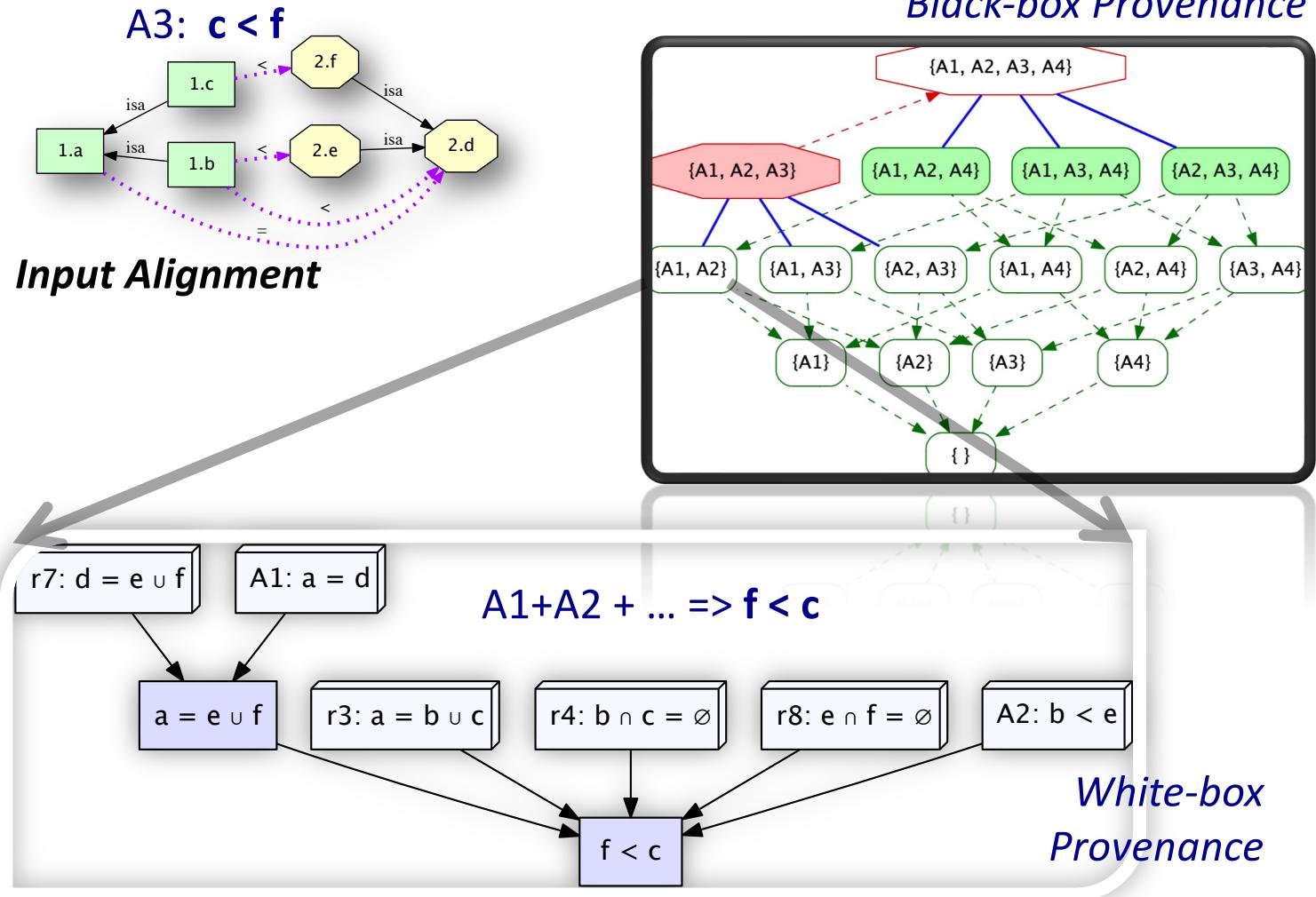
Hybrid Provenance



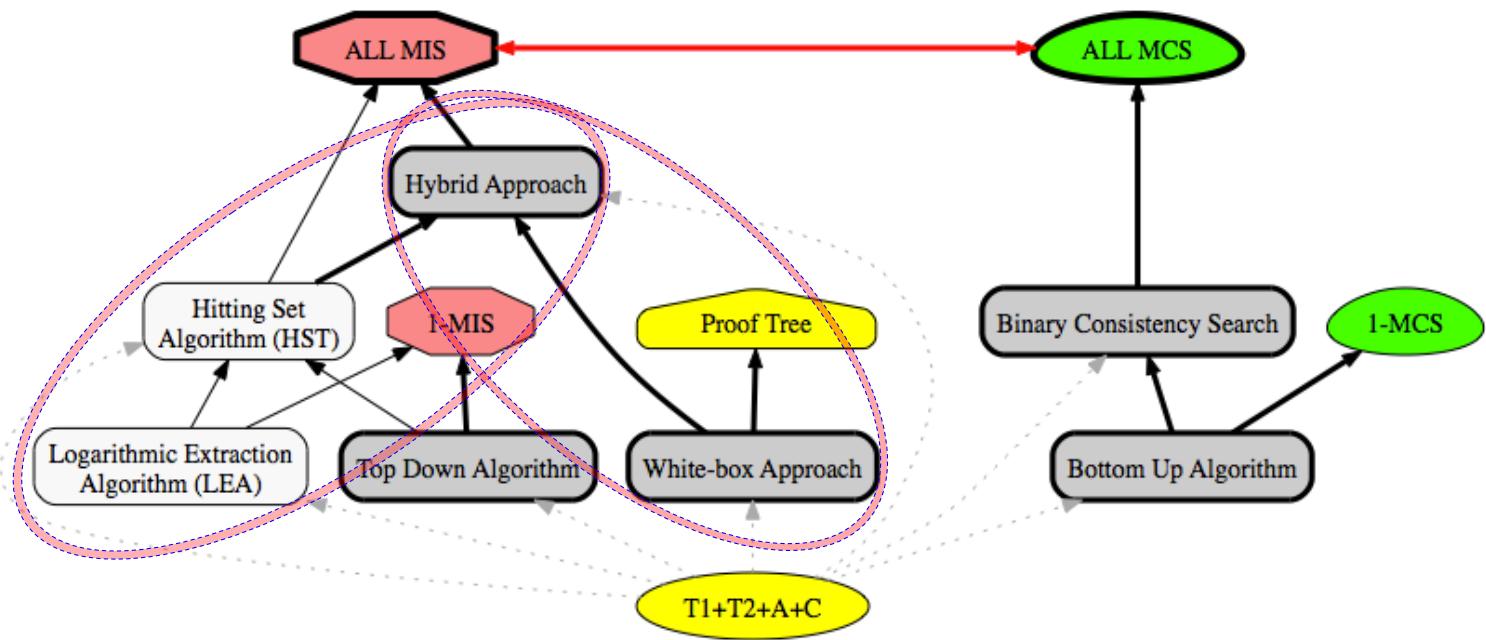
Black-box Provenance



Hybrid Provenance



The Hybrid Approach



Hybrid Approach

Algorithm 3 White-Box Approach

Input: System description SD , a set of constraints C

Output: All diagnosis proof trees

`ComputeAllProofTrees(SD, C):`

- 1: Encode SD and C in Datalog rules
- 2: Rewrite Datalog rules to ones with provenance
- 3: Run ASP reasoner to get boolean expressions for NOK
- 4: Construct diagnosis proof trees using the boolean expressions

Algorithm 4 Hybrid Approach

Input: System description SD , a set of constraints C

Output: All diagnoses (MIS)

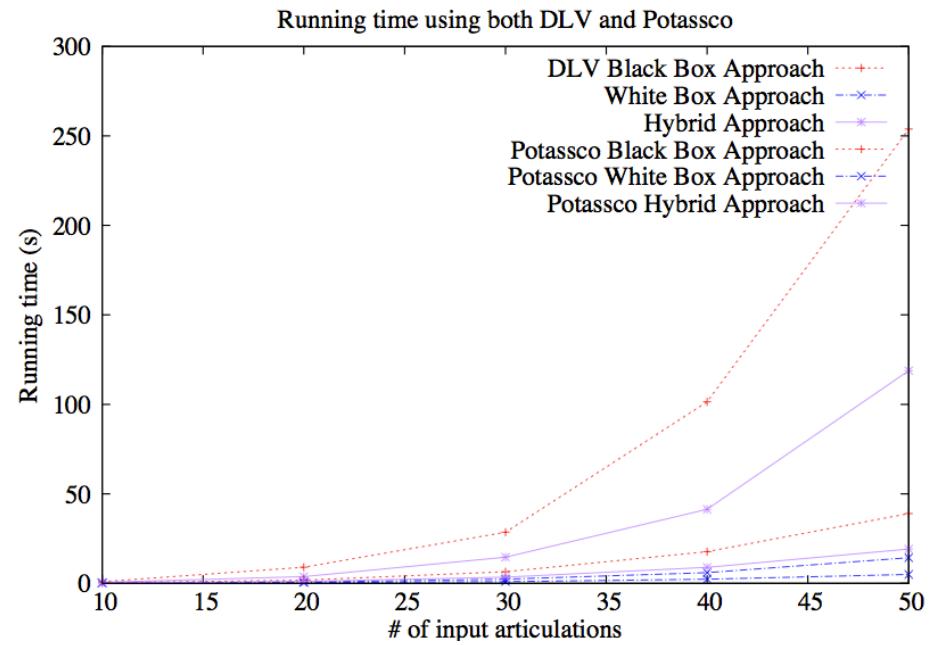
`ComputeAllMISHybrid(SD, C):`

- 1: $T_s \leftarrow \text{ComputeAllProofTrees}(SD, C)$
- 2: $C' \leftarrow$ set of leaf nodes of the proof trees T_s
- 3: **return** $\text{ComputeAllMIS}(SD, C')$

What articulations
contribute to some
inconsistency?

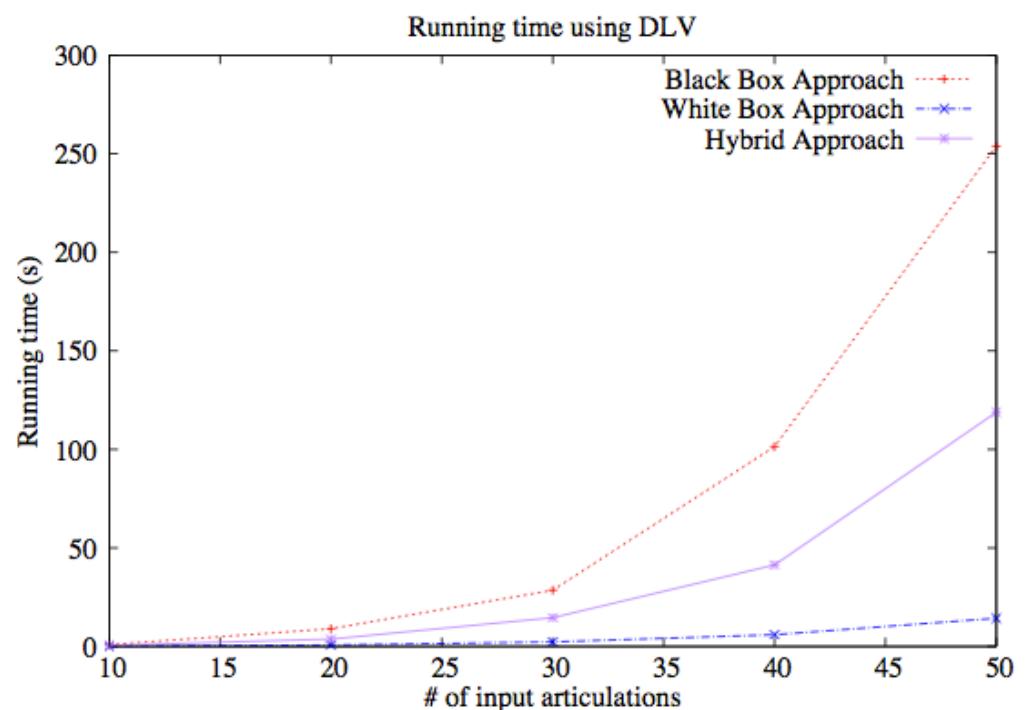
Good old black-box
(HST)

Benchmark Results



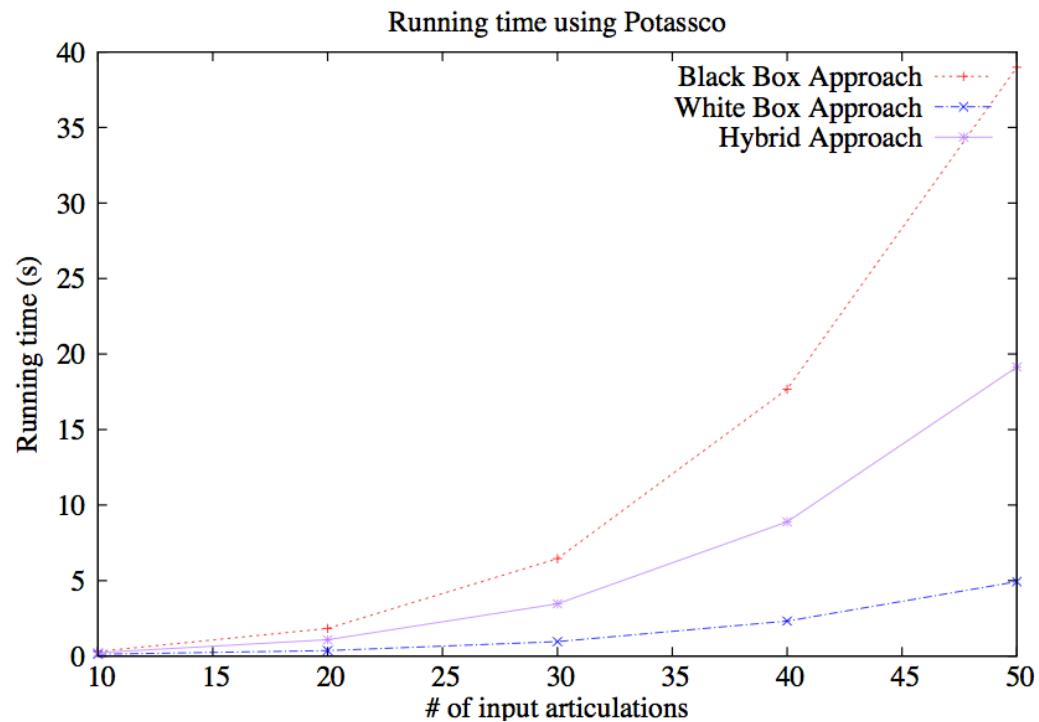
- White-box < **Hybrid** < **Black-box**
(runtimes)
- Note: white-box does **not** give you a diagnosis
- Potassco < DLV

Benchmark DLV



- White-box < Hybrid < Black-box (runtimes)
- Potassco < DLV

Benchmark Clingo



- White-box < Hybrid < Black-box (runtimes)
- Potassco < DLV

Conclusions

- ASP rules can be used to efficiently solve real-world taxonomy reasoning problems
- Reiter's diagnosis useful to debug inconsistent alignments
- **Adding a “white-box” provenance approach speeds up state-of-the-art HST algorithm by eliminating independent articulations**
- Future work:
 - Further improvements, including **parallelism**:
 - Trade-off with **sharing** inferences across parallel instances