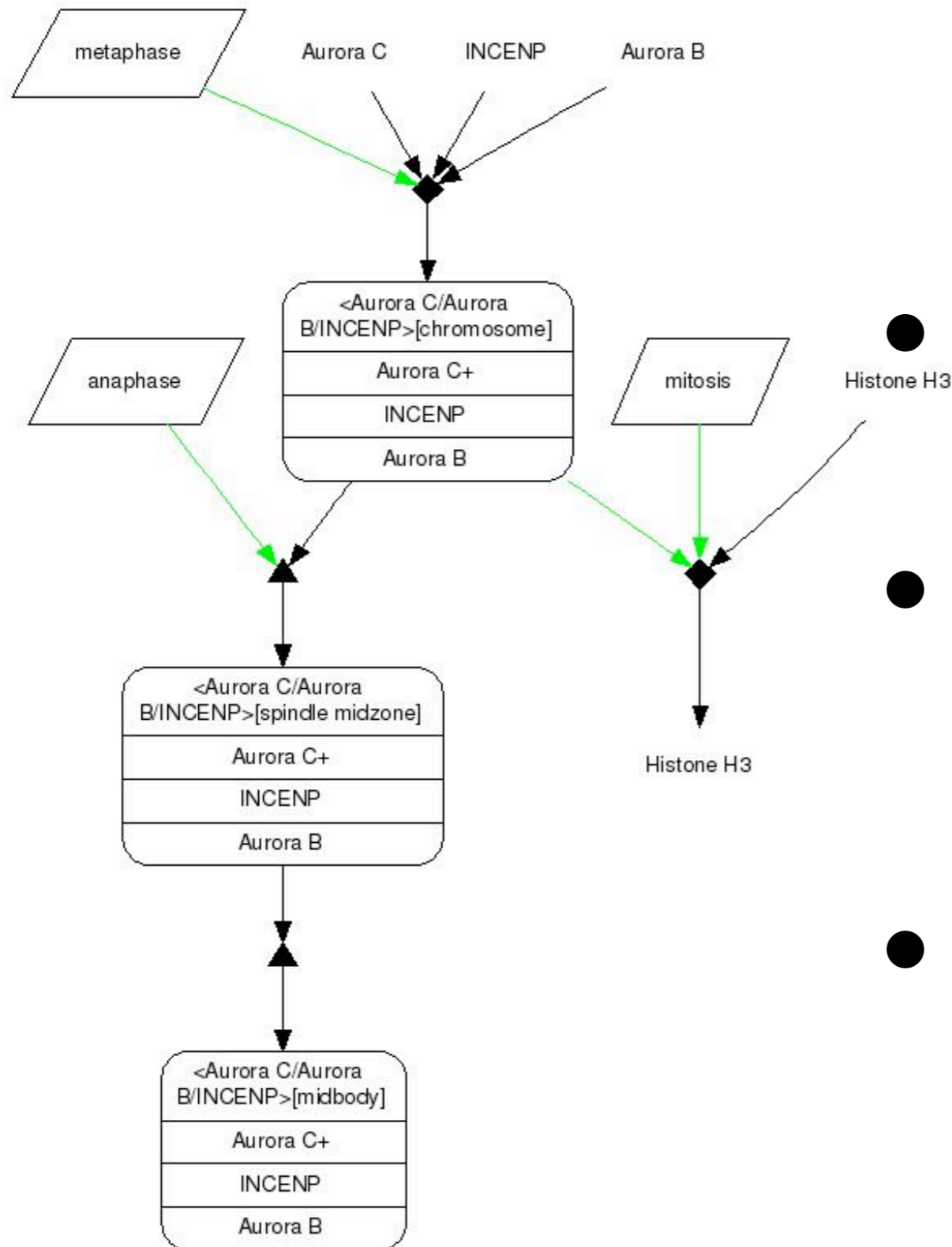# An attempt to use literature curated pathway DBs

Charles Vaske
Stuart Lab Meeting
May 6th, 2009

# Structured Pathways
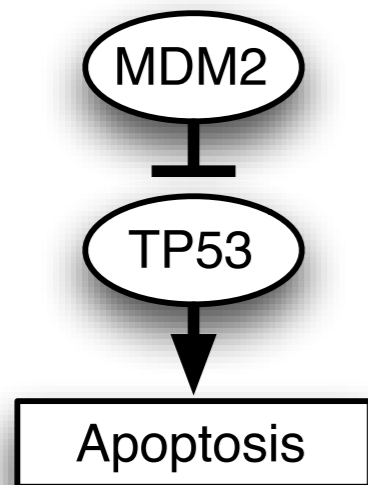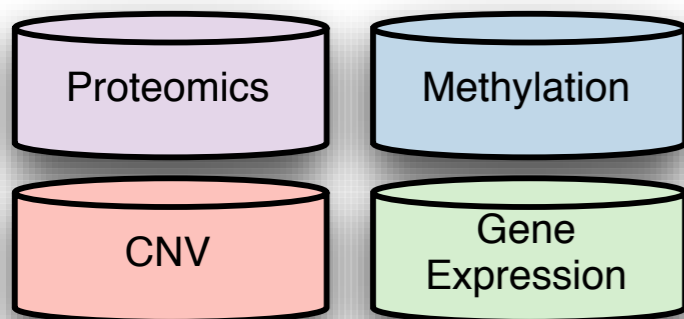


- Lots of cancer research/ genes/data

- Subsequently, we know a lot about pathways active in cancer

- Can we use this structured knowledge?
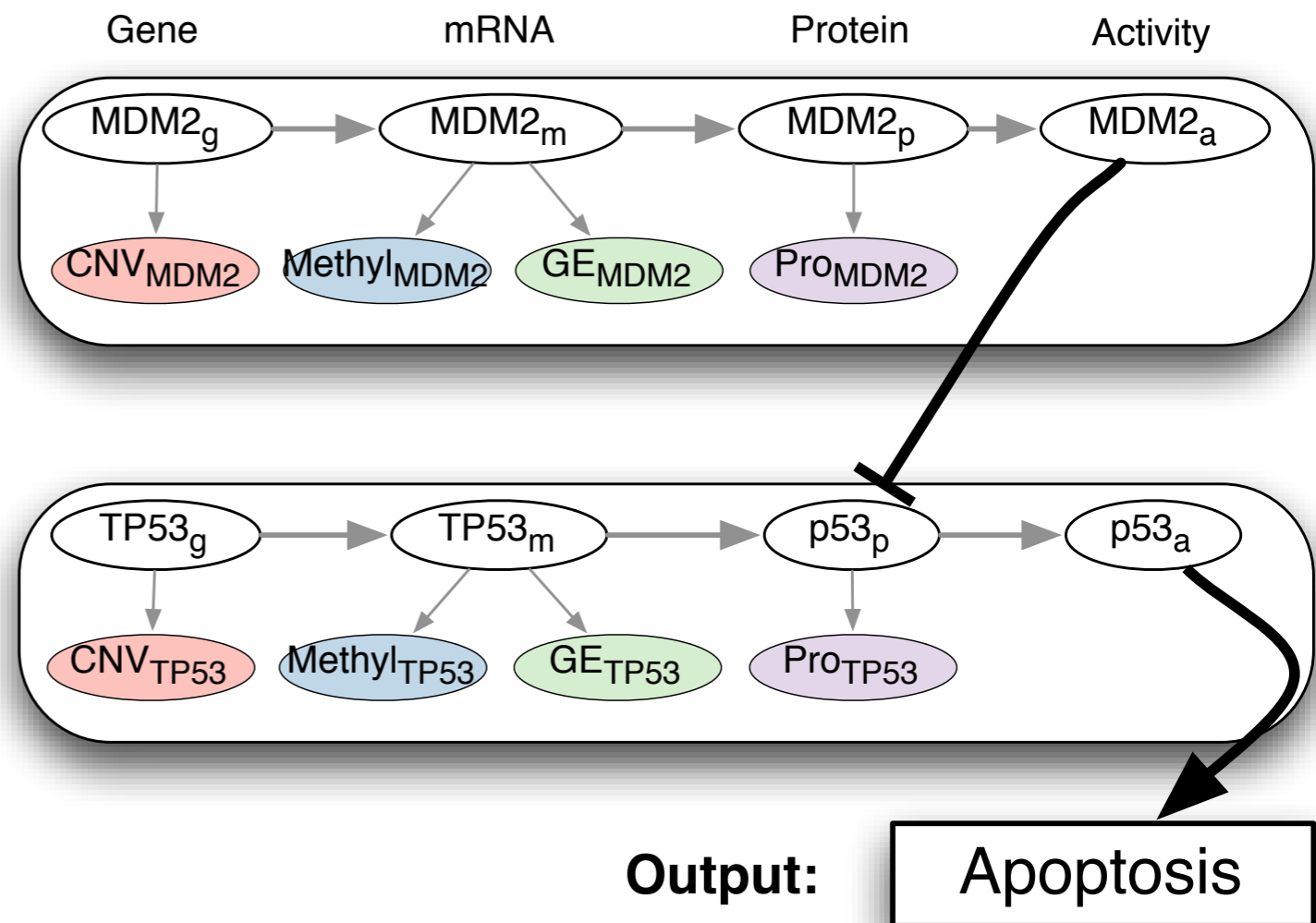
# Use in clinical samples

# Outline

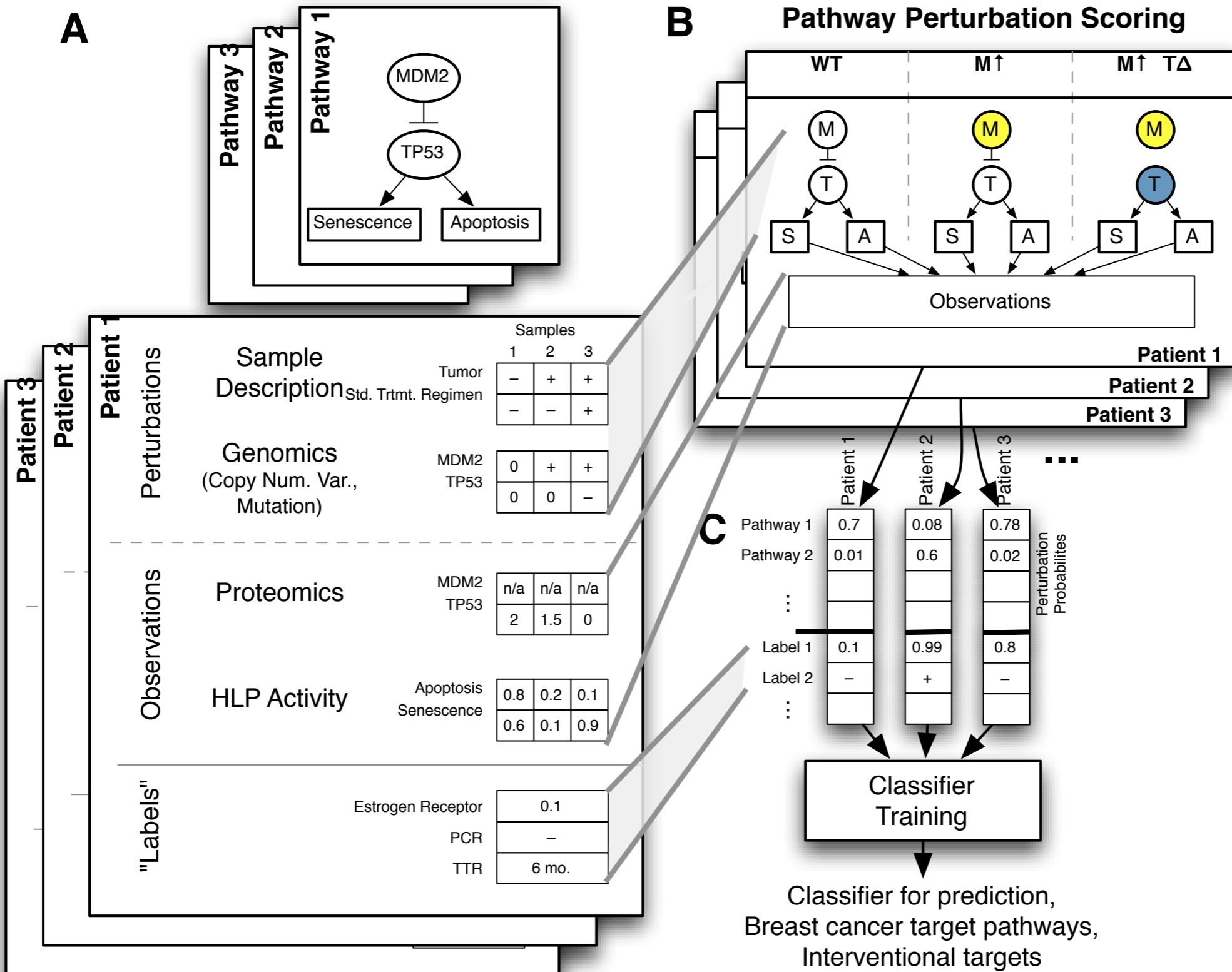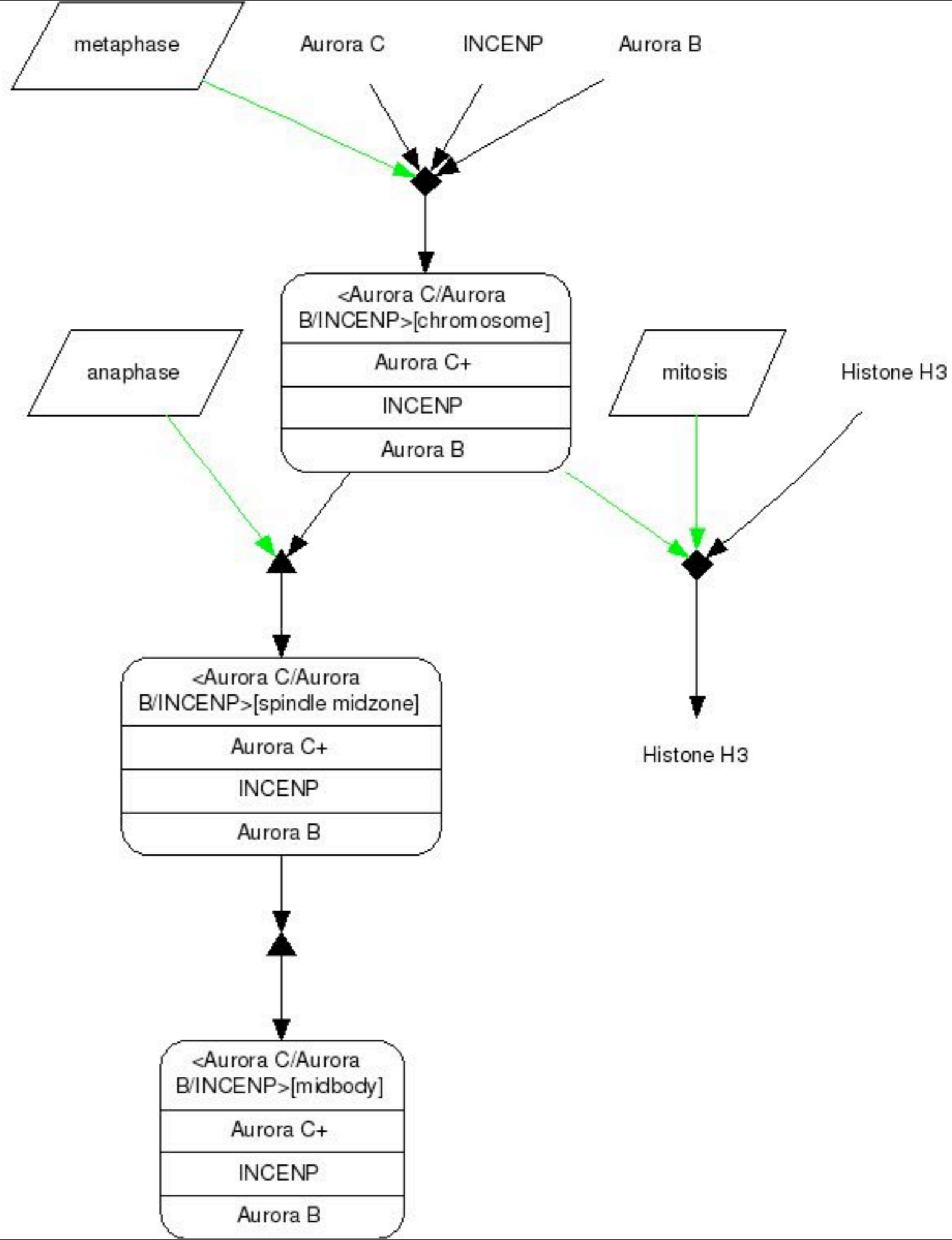1. Get pathways (ugly, 50%-95% done)

2. Convert to graphical model

3. Add evidence from patient

4. Infer the value of hidden variables (i.e. Apoptosis, Chemotaxis)

5. Solve cancer (finally)

# BioPAX

```
<?xml version="1.0" encoding="UTF-8"?>
<rdf:RDF xmlns:rdf="http://www.w3.org/1999/02/22-rdf-syntax-ns#"
  xmlns:bp="http://www.biopax.org/release/biopax-level2.owl#"
  xmlns:rdfs="http://www.w3.org/2000/01/rdf-schema#"
  xmlns:owl="http://www.w3.org/2002/07/owl#"
  xmlns:xsd="http://www.w3.org/2001/XMLSchema#"
  xmlns="http://pid.nci.nih.gov/biopax#"
  xml:base="http://pid.nci.nih.gov/biopax">
  <owl:Ontology rdf:about="">
    <owl:imports rdf:resource="http://www.biopax.org/release/biopax-level2.owl" />
    <rdfs:comment rdf:datatype="http://www.w3.org/2001/XMLSchema#string">BioPAX output created 2009_04_14 11:44::00, converted from the Pathway Interaction Database, National Cancer Institute, http://pid.nci.nih.gov.</rdfs:comment>
  </owl:Ontology>
  <bp:bioSource rdf:ID="Homo_sapiens">
    <bp:NAME rdf:datatype="http://www.w3.org/2001/XMLSchema#string">Homo sapiens</bp:NAME>
    <bp:TAXON-XREF rdf:resource="#NCBI_taxonomy_9606" />
  </bp:bioSource>
  <bp:unificationXref rdf:ID="NCBI_taxonomy_9606">
    <bp:DB rdf:datatype="http://www.w3.org/2001/XMLSchema#string">NCBI_taxonomy</bp:DB>
    <bp:ID rdf:datatype="http://www.w3.org/2001/XMLSchema#string">9606</bp:ID>
  </bp:unificationXref>
  <bp:dataSource rdf:ID="PID_DataSource">
    <bp:NAME rdf:datatype="http://www.w3.org/2001/XMLSchema#string">Pathway Interaction Database</bp:NAME>
    <bp:COMMENT rdf:datatype="http://www.w3.org/2001/XMLSchema#string">http://pid.nci.nih.gov</bp:COMMENT>
  </bp:dataSource>
  <bp:dataSource rdf:ID="PID_Curated_DataSource">
    <bp:NAME rdf:datatype="http://www.w3.org/2001/XMLSchema#string">Pathway Interaction Database NCI-Nature Curated Data</bp:NAME>
    <bp:COMMENT rdf:datatype="http://www.w3.org/2001/XMLSchema#string">http://pid.nci.nih.gov</bp:COMMENT>
  </bp:dataSource>
  <bp:dataSource rdf:ID="PID_BioCarta_DataSource">
    <bp:NAME rdf:datatype="http://www.w3.org/2001/XMLSchema#string">Pathway Interaction Database BioCarta Data</bp:NAME>
    <bp:COMMENT rdf:datatype="http://www.w3.org/2001/XMLSchema#string">
```
-u:--  ac.owl        Top L1      (nXML Valid)----------------------------------
Using schema ~/.emacs.lisp/nxml-mode-20041004/schema/rdfxml.rnc
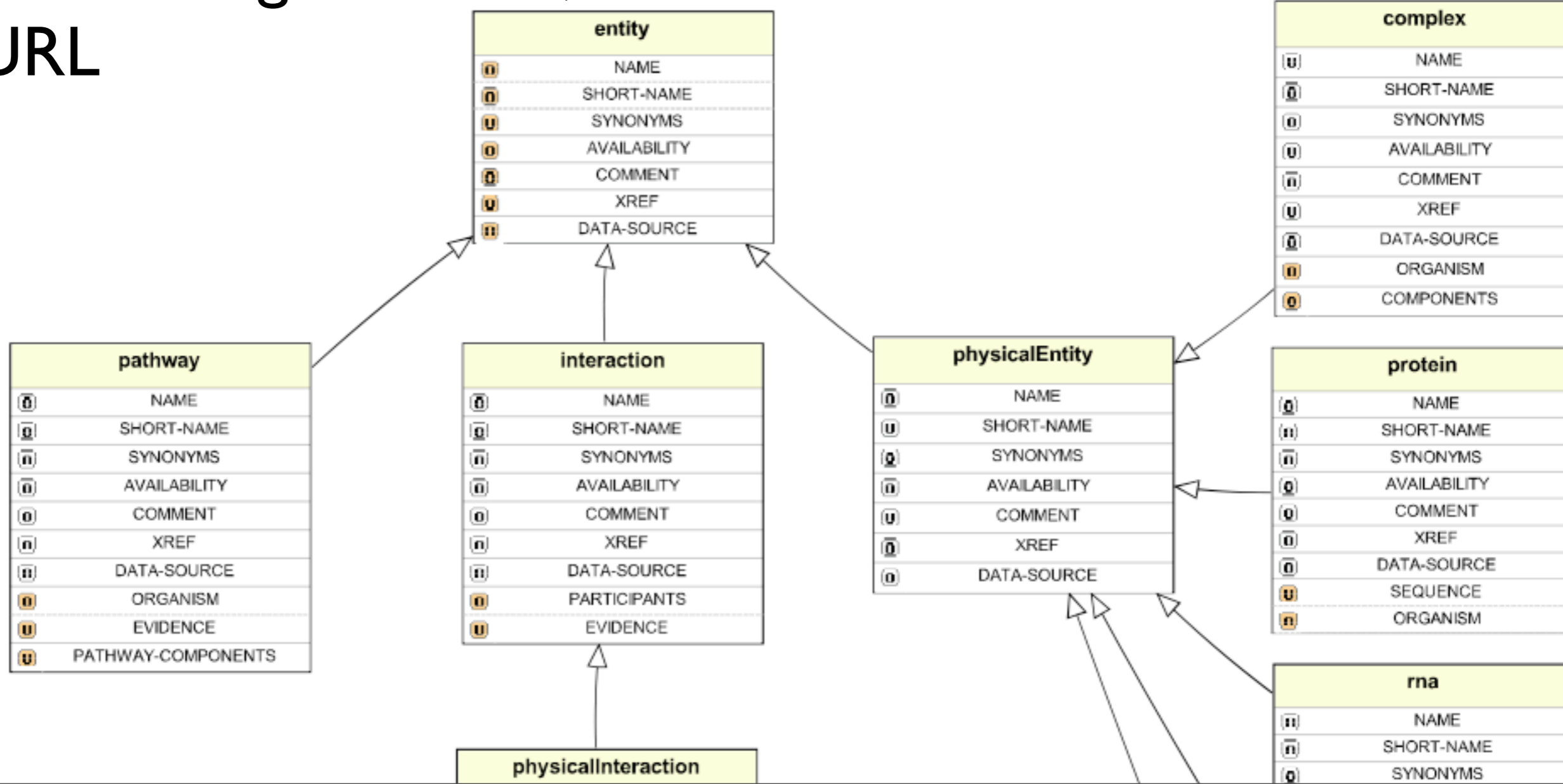
Emacs@dhcp-63-190.cse.ucsc.edu

- Based on OWL
  *Web Ontology Lang.*

- Based on RDF
  *Resource Desc. Format*

- Not human-readable

- Must use tools!

- I love to complain about it

- Three levels (versions), people only use level 2 (I think)

- Defines "things" which have various properties, including a "class"

- Each "thing" is a URI, which looks like a URL
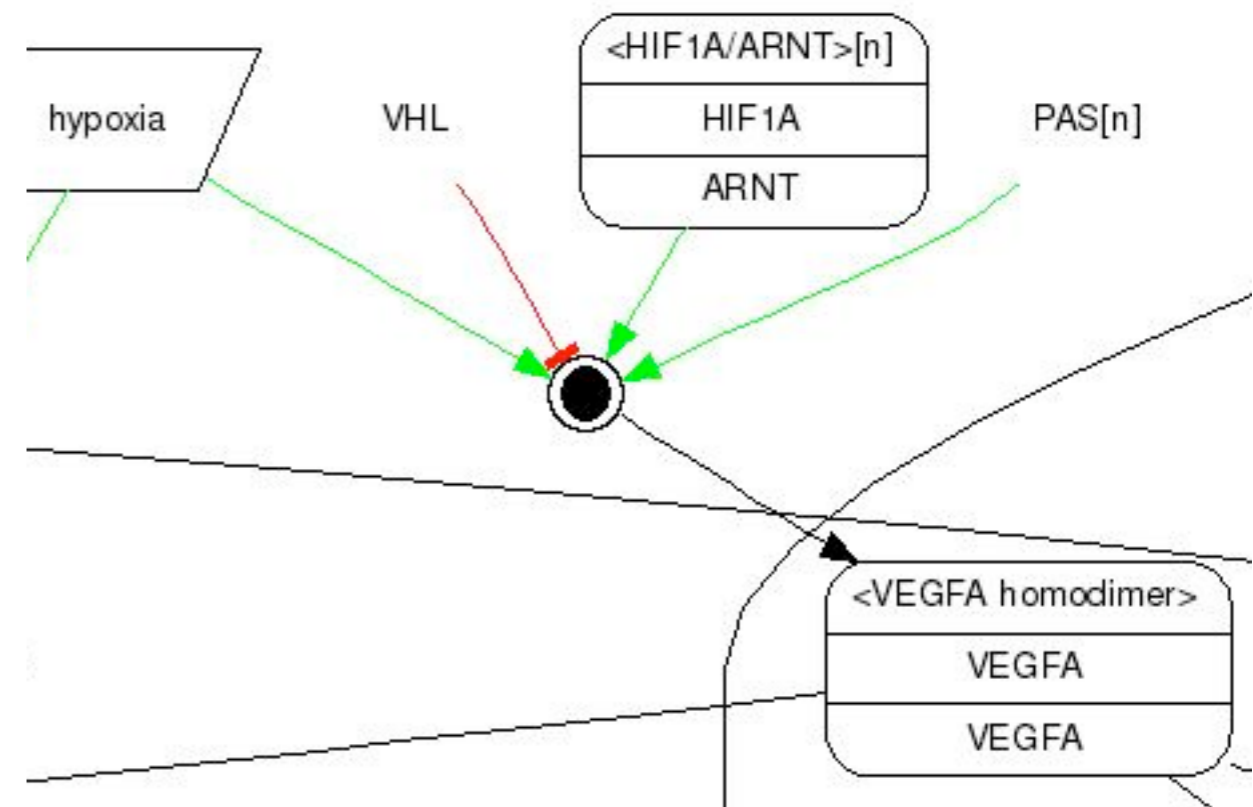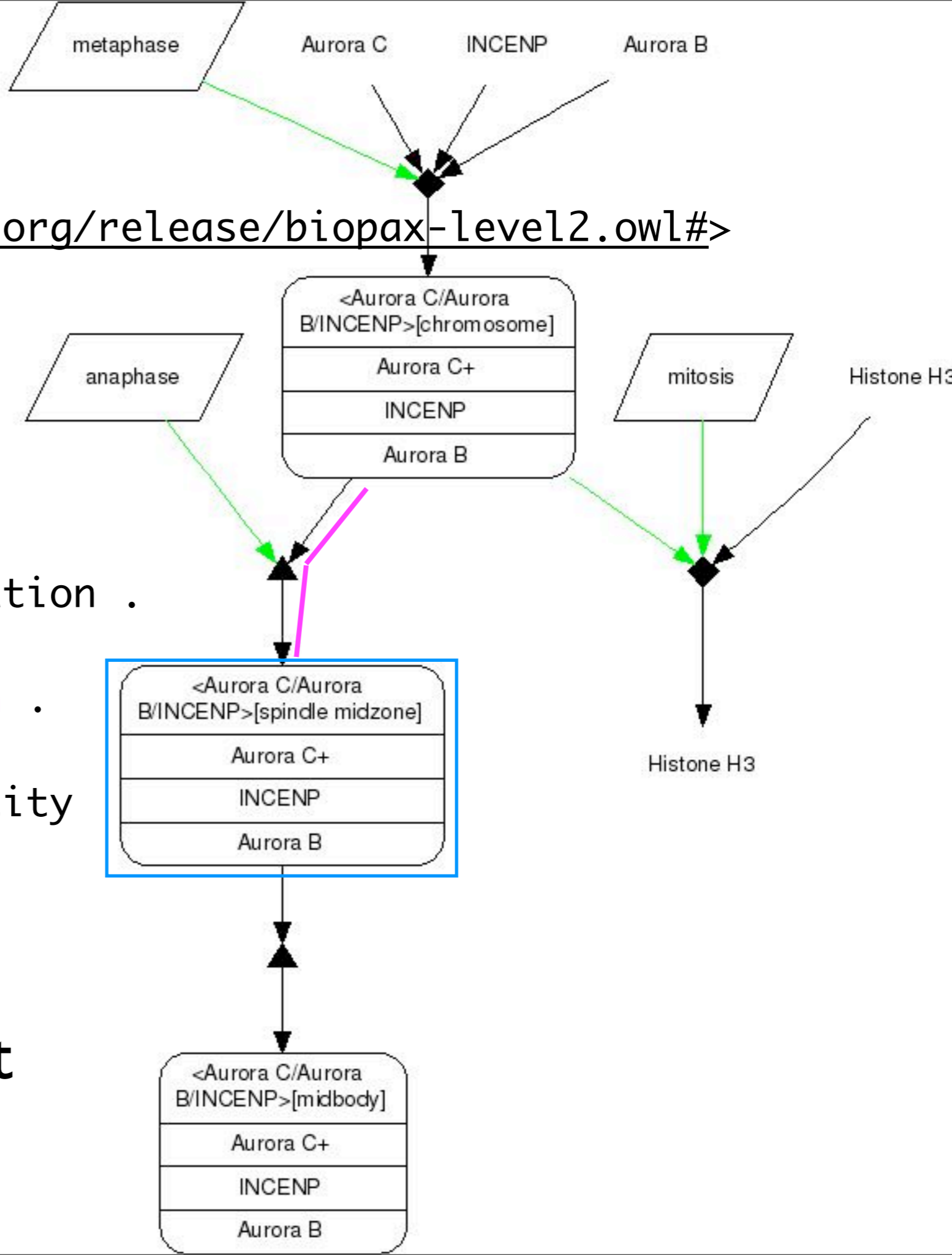
# BioPAX

# RDF/OWL/BioPAX Tools

- **Protege**: from Stanford, designed more for creating a BioPAX more than looking at "data" in that "format"

- **SPARQL/roqet**: sort of like SQL for RDF. Don't use XML tools, you may miss things due to variations in serializations.

# Caveat

- All this dense typing and formating is *extremely* expressive

- However, the amount of expression impedes programmatic understanding

- Test, test, test



This shows the "transcription" of a complex. The meaning is obvious to a human, but befuddling to my naive scripts.

```
PREFIX bp: <http://www.biopax.org/release/biopax-level2.owl#>

SELECT
    ?goname
    ?entity
    ?activation
WHERE {
    ?mod bp:CONTROL-TYPE ?activation .
    ?mod bp:NAME ?goname .
    ?mod bp:CONTROLLED ?reaction .
    ?reaction bp:RIGHT ?pep .
    ?pep bp:PHYSICAL-ENTITY ?entity
}
```

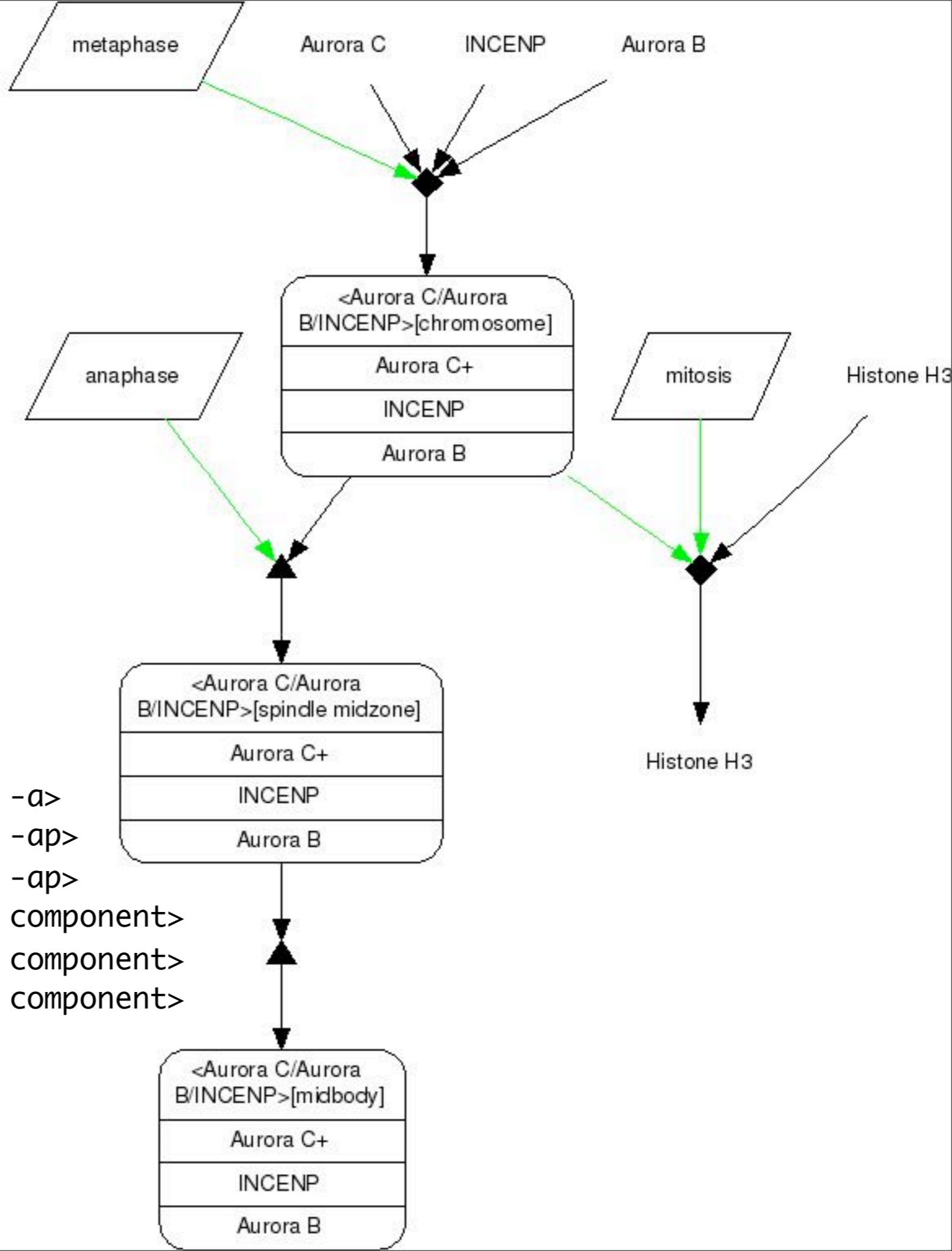Example query: find abstract processes that are parents

# Parsing

- Started by finding the proper queries to extract interactions, names, parts of complexes...

- Want a simple tab-delimited format:

| | | |
|---|---|---|
| abstract | metaphase | |
| abstract | mitosis | |
| complex | AurC/AurB/INCENP | |
| protein | H3F3A | |
| protein | AuroraB | |
| protein | AuroraC | |
| protein | INCENP | |

**Entity Definitions**

| | | |
|---|---|---|
| AurC/AurB/INCENP | H3F3A | -a> |
| mitosis | H3F3A | -ap> |
| metaphase | AurC/AurB/INCENP | -ap> |
| AuroraB | AurC/AurB/INCENP | component> |
| INCENP | AurC/AurB/INCENP | component> |
| AuroraC | AurC/AurB/INCENP | component> |

**Entity Interactions**

PID Aurora C signaling

| abstract | metaphase | |
| abstract | mitosis | |
| complex | AurC/AurB/INCENP | |
| protein | H3F3A | |
| protein | AuroraB | |
| protein | AuroraC | |
| protein | INCENP | |
| AurC/AurB/INCENP | H3F3A | -a> |
| mitosis | H3F3A | -ap> |
| metaphase | AurC/AurB/INCENP | -ap> |
| AuroraB | AurC/AurB/INCENP | component> |
| INCENP | AurC/AurB/INCENP | component> |
| AuroraC | AurC/AurB/INCENP | component> |

# My hopeful monster



- Makefile converted to executable script

- A bit experimental

# $MAPDIR/Data/Pathways

- Early, molten stage, but useful

- Human/NCIPID has NCI pathways
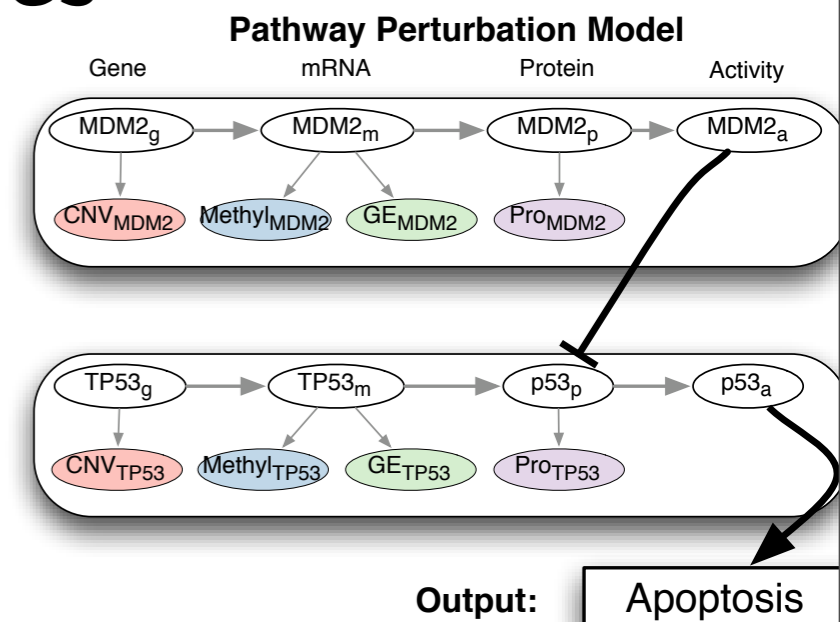
- Human/KEGG has early KEGG attempts

# Pathway stats



| | |
|---|---|
| 132 | **Pathways** |
| 3766 | **Unique Entities** |
| 7569 | **Entity instances** |
| 10182 | **Unique Interactions** |
| | **Entity Breakdown** |
| 1742 | protein |
| 1638 | complex |
| 296 | abstract |
| 90 | chemical |
| | **Interaction Breakdown** |
| 2619 | -a> (activation) |
| 278 | -a| (inhibition) |
| 874 | -ap> (abstract) |
| 103 | -ap| |
| 528 | -t> (transcription) |
| 104 | -t| |
| 5676 | component> |

Number of interactions in pathway

Number of entities in pathway

# Outline

~~1. Get pathways (ugly, 50%-95% done)~~

2. Convert to graphical model

3. Add evidence from patient
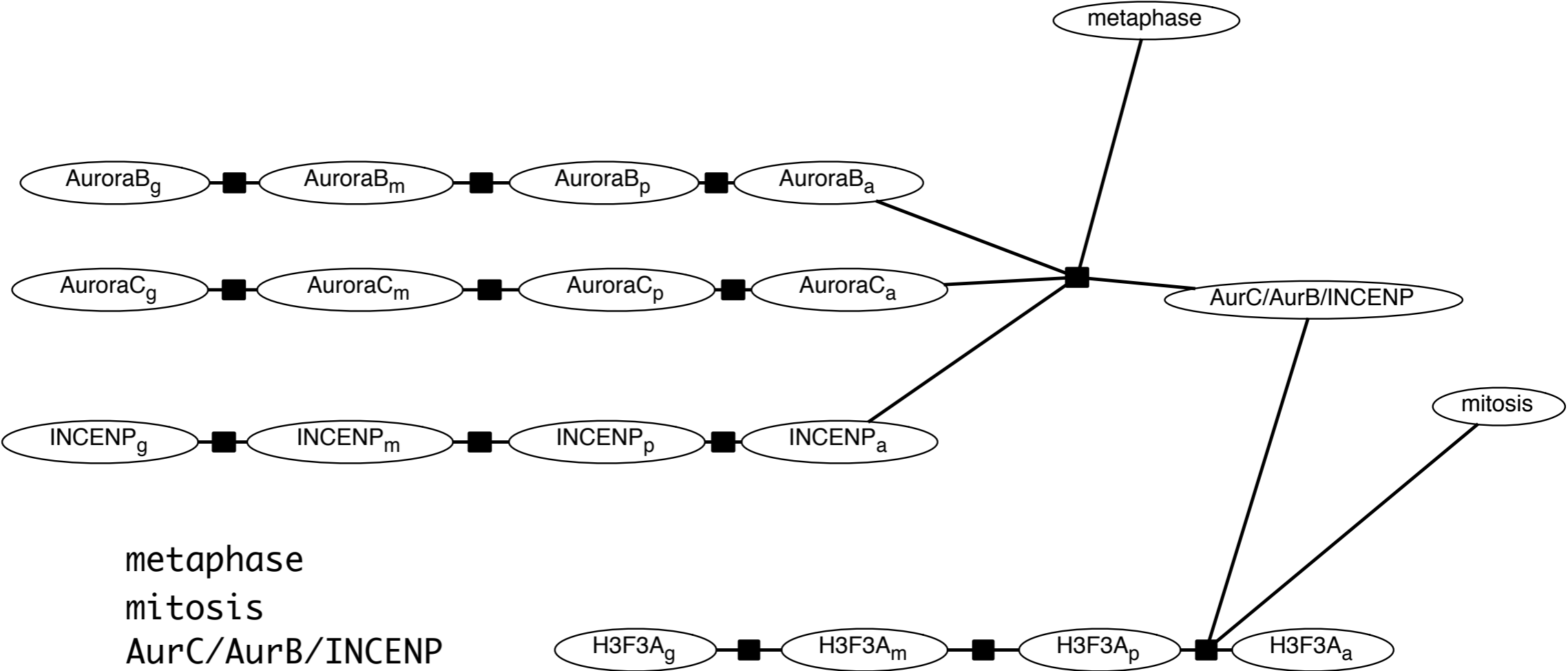
4. Infer the value of hidden variables (i.e. Apoptosis, Chemotaxis)

5. Solve cancer (finally)

# Aurora C Factor Graph



| | | |
|---|---|---|
| abstract | metaphase | |
| abstract | mitosis | |
| complex | AurC/AurB/INCENP | |
| protein | H3F3A | |
| protein | AuroraB | |
| protein | AuroraC | |
| protein | INCENP | |
| AurC/AurB/INCENP | H3F3A | -a> |
| mitosis | H3F3A | -ap> |
| metaphase | AurC/AurB/INCENP | -ap> |
| AuroraB | AurC/AurB/INCENP | component> |
| INCENP | AurC/AurB/INCENP | component> |
| AuroraC | AurC/AurB/INCENP | component> |

# Aurora C Evidence



| | | |
|---|---|---|
| AuroraB | genome | 0.87082 |
| AuroraB | mRNA | 0.37673 |
| AuroraC | genome | 0.170729 |
| AuroraC | mRNA | 0.045578 |
| INCENP | genome | -0.082277 |
| INCENP | mRNA | -0.060272 |
| H3F3A | genome | -0.411328 |

- Data points are signed, log p-values

- Right now, I discretize into up/down/same at 0.05 level

- Therefore, many patients look "identical" on hidden variables

# Aurora C Inference

- using the package libDAI, which implements many approximate inference algorithms (and exact)

- Using exact at the moment

- 128 patients, 132 pathways ~ 2 hours

# Prelim. Pathway results

- 2 data sets

  - Glioblastoma        224 samples

  - Ovarian Cancer    128 samples

- Still working out kinks in pipeline

- Not satisfied with data treatment