

chado












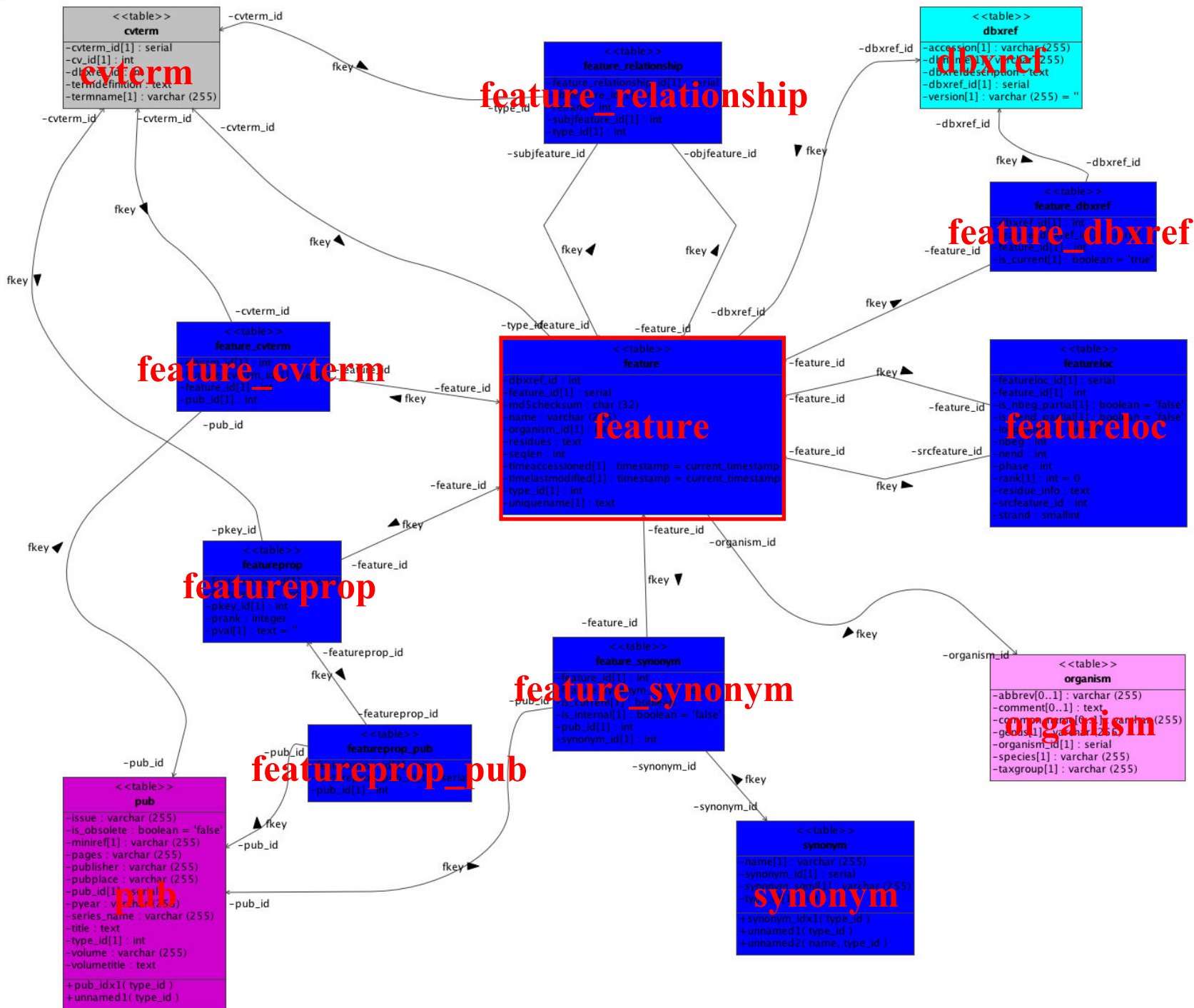
Generic model organism database schema

**presented at Cold Spring Harbor GMOD meeting, 5/03
by Stan Letovsky, sletovsky@aol.com**

Chado modules



 Genetics Module	Alleles, and relationships between alleles and phenotypes
 Sequence Module	Anything related to biological sequences and annotation
 Map Module	Any kind of localisation excluding sequence localisations
 Expression Module	Transcription events, including space/time localisation. Also for protein expression
 Companalysis Module	Adjunct to sequence module for in-silico analysis
 CV Module	Controlled Vocabularies / Ontologies
 Organism Module	Data related to taxonomy / species
 Pub Module	Publication / Biblio / References
 General Module	General / Core





Sequence_Feature_Ontology

- i** chromosome variation
- i** consequences of mutation
- i** sequence feature
 - i** located sequence feature
 - i** annotation comment
 - i** gene element
 - i** regulatory region
 - i** transcript region
 - i** nucleotide motif
 - i** reagent
 - i** region
 - i** conserved region
 - i** deletion
 - i** exon
 - i** gene group
 - i** gene region
 - i** integrated virus
 - i** intron
 - i** match
 - i** non-gene feature
 - i** recombination feature
 - i** repeat region
 - i** replication feature
 - i** sequence secondary structure
 - i** sequence variation feature
 - i** site
 - i** sequence attribute

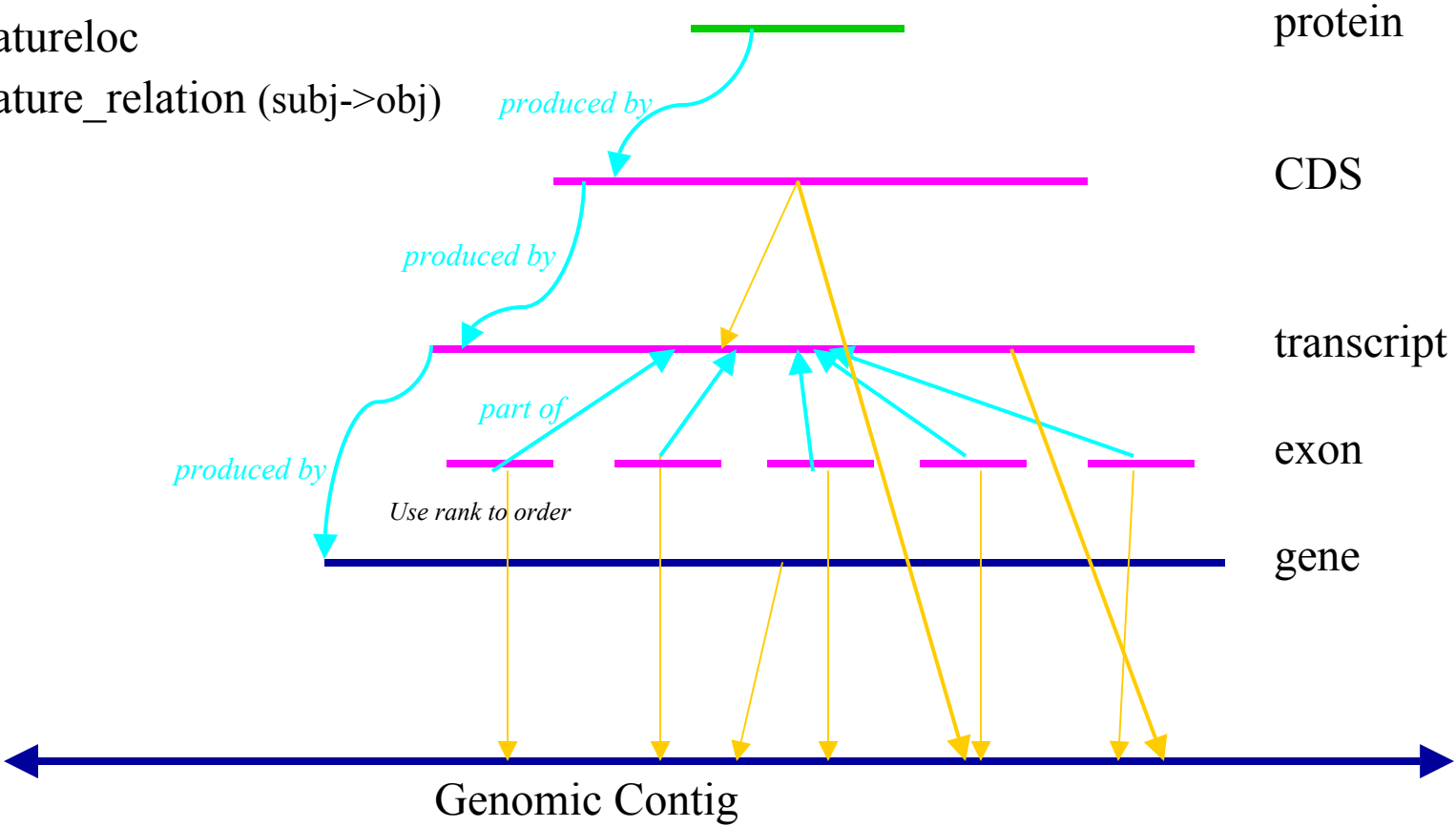
Sequence Ontology

- i** region
 - i** conserved region
 - i** coding conserved region
 - i** noncoding conserved region
 - i** syntenic region
 - i** deletion
 - i** exon
 - i** 3'-exon
 - i** 5'-exon
 - i** coding exon
 - i** interior exon
 - i** noncoding exon
 - i** single exon
 - i** gene group
 - i** gene region
 - i** integrated virus
 - i** intron
 - i** 3'-intron
 - i** 5'-intron
 - i** autocatalytically spliced intron
 - P** branch site
 - i** group I intron
 - i** group II intron
 - i** interior intron
 - i** intronic splice enhancer
 - P** polypyrimidine tract
 - i** match
 - i** nucleotide to nucleotide match
 - i** nucleotide to protein match
 - i** protein to nucleotide match
 - i** protein to protein match



Central Dogma: single spliced transcript

- Feature (Colors: DNA, RNA, Protein)
- Featureloc
- Feature_relation (subj->obj)

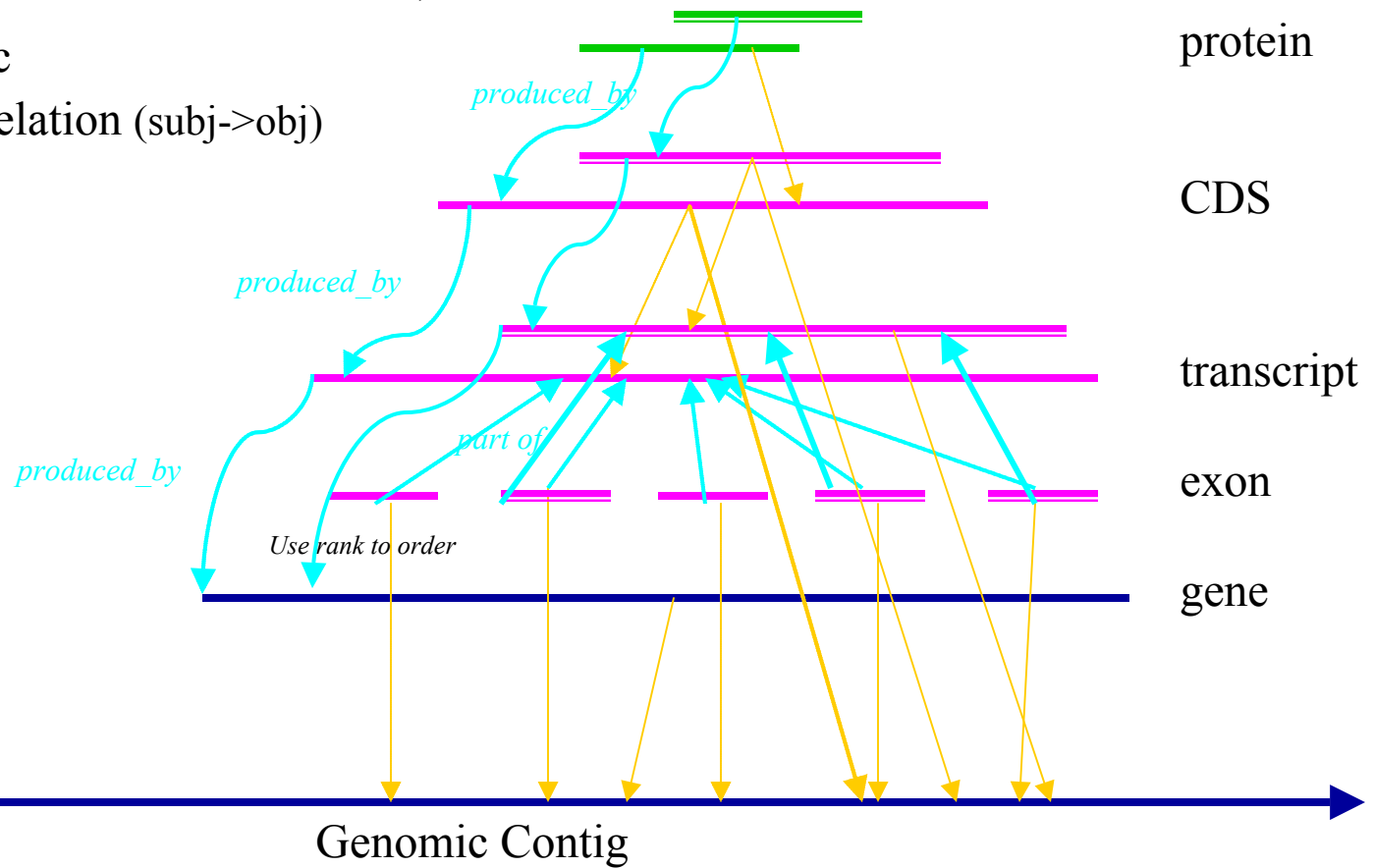


Central Dogma: 2nd transcript (alt. Splicing)

— Feature (Colors: DNA, RNA, Protein)

→ Featureloc

→ Feature_relation (subj->obj)

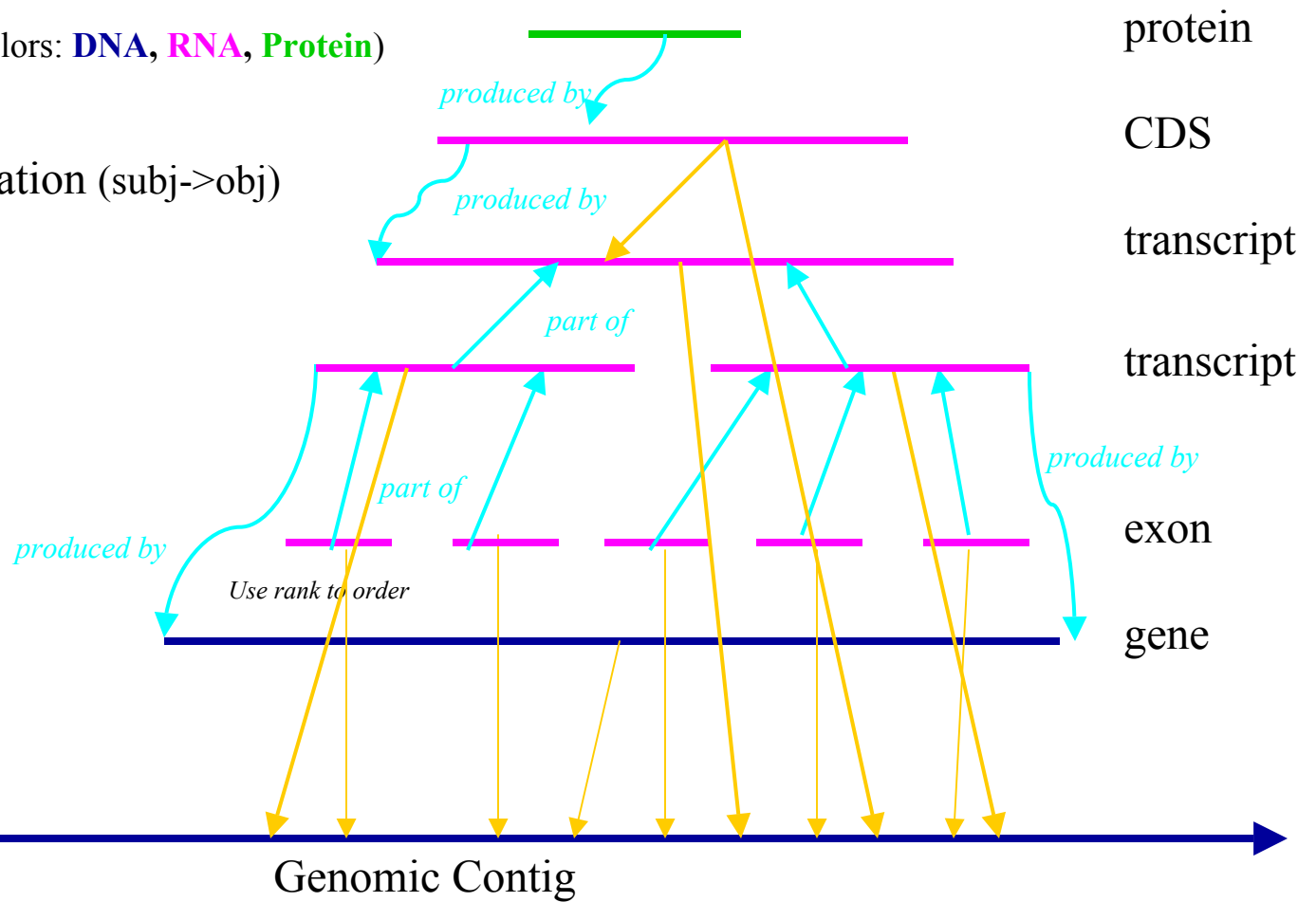


Pathological Cases: trans-splicing

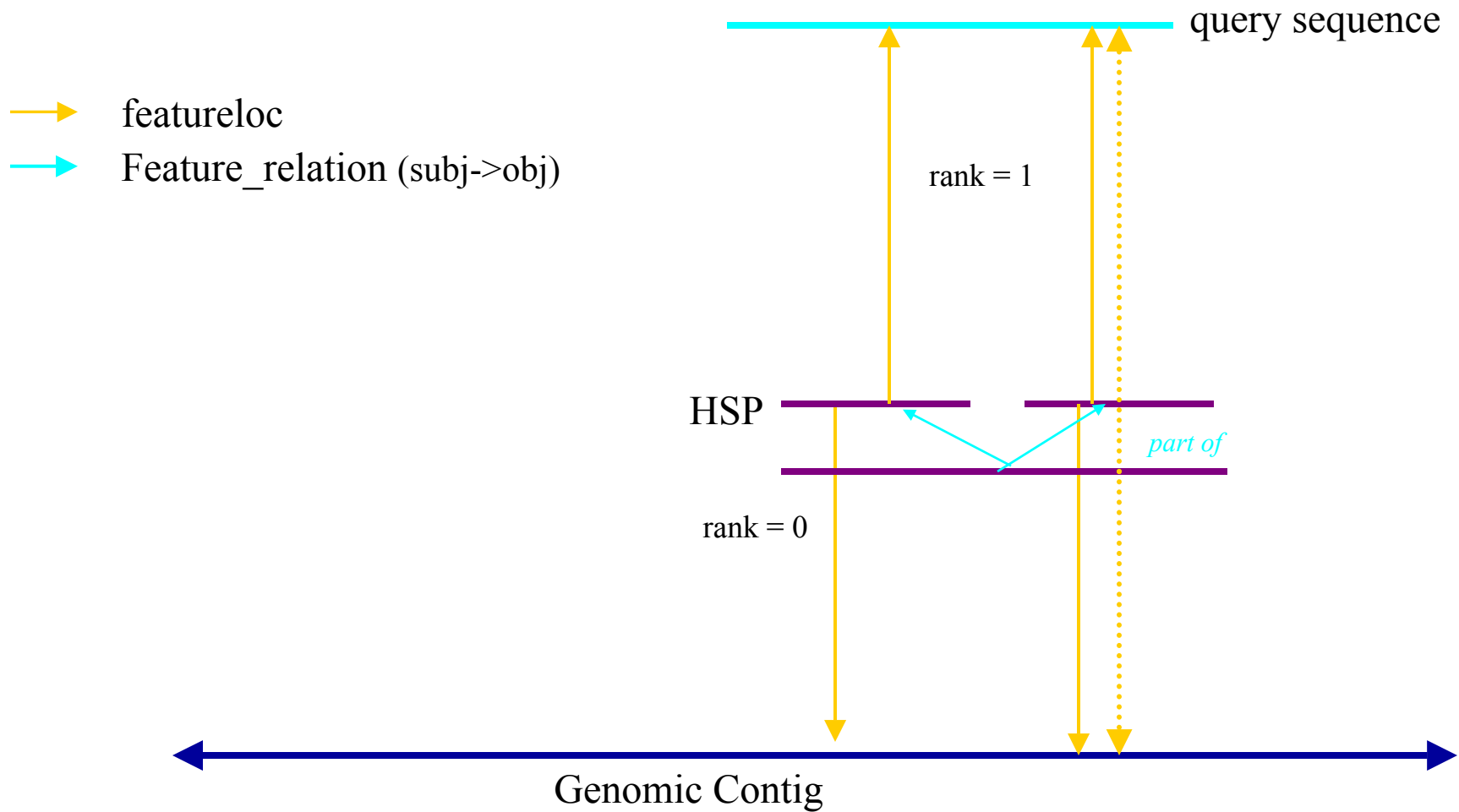
— Feature (Colors: DNA, RNA, Protein)

→ Featureloc

→ Feature_relation (subj->obj)

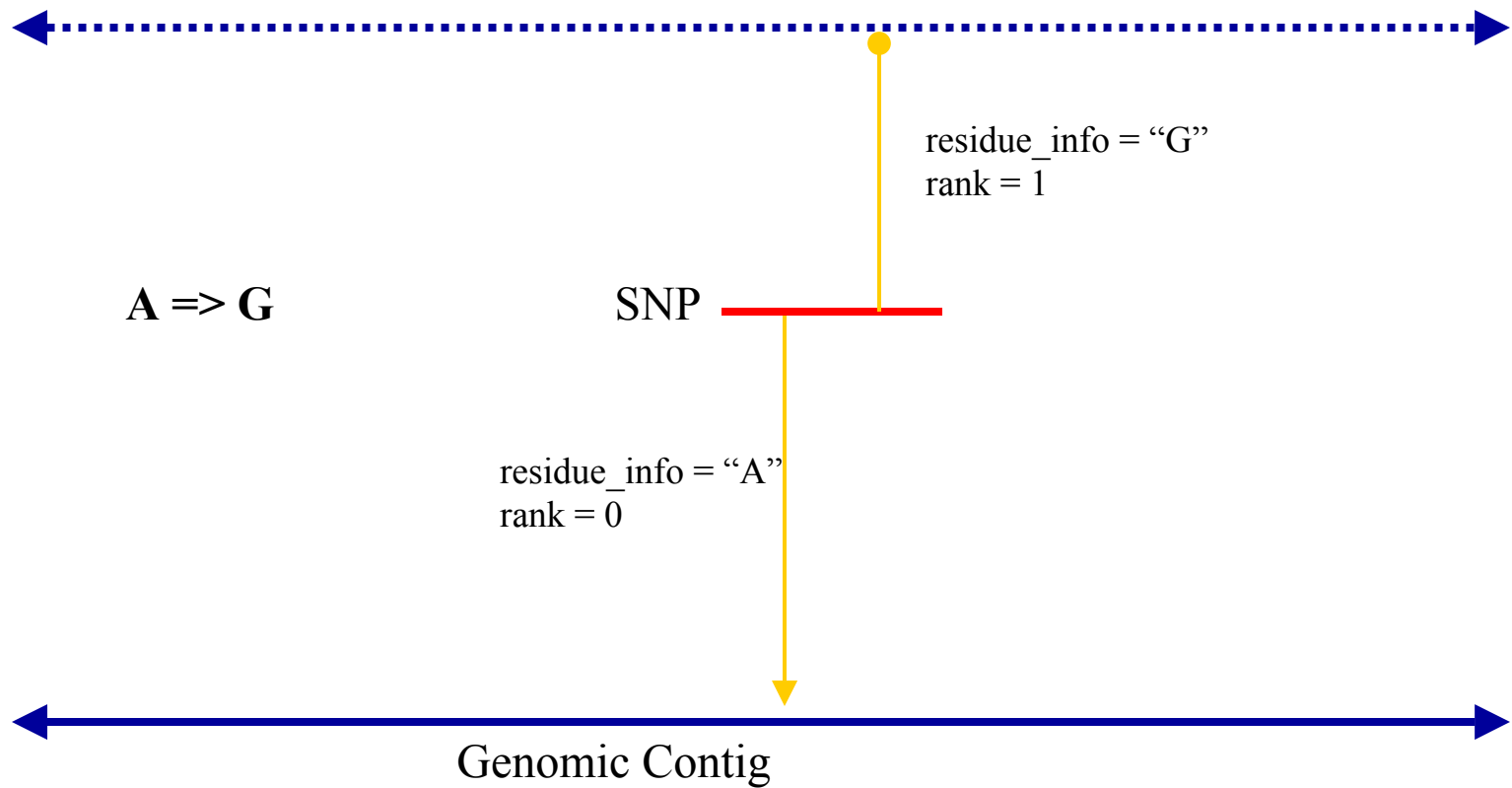


Pairwise Alignments



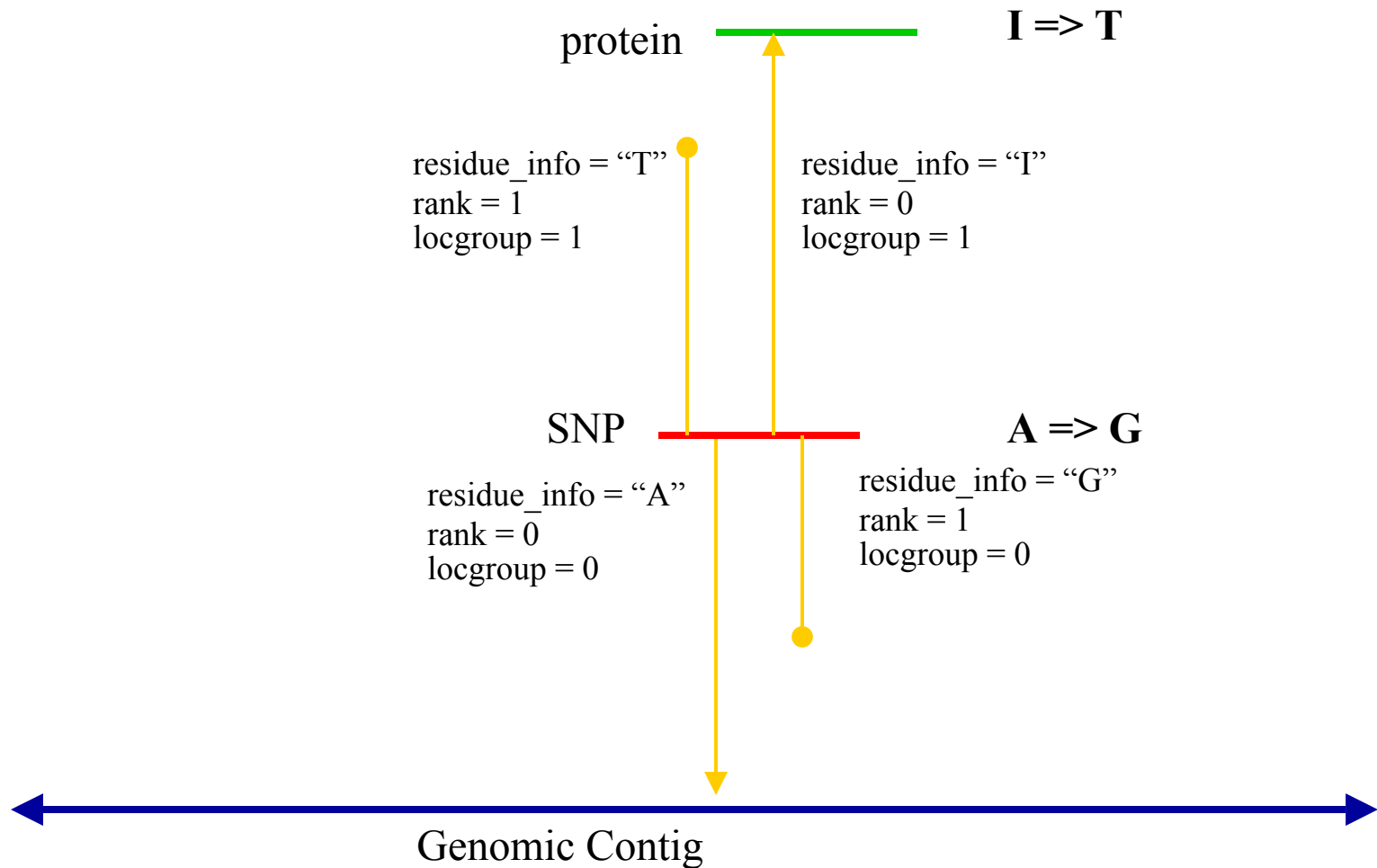
Sequence variations

SNPs



Sequence variations

SNPs (redundant mapping to protein)



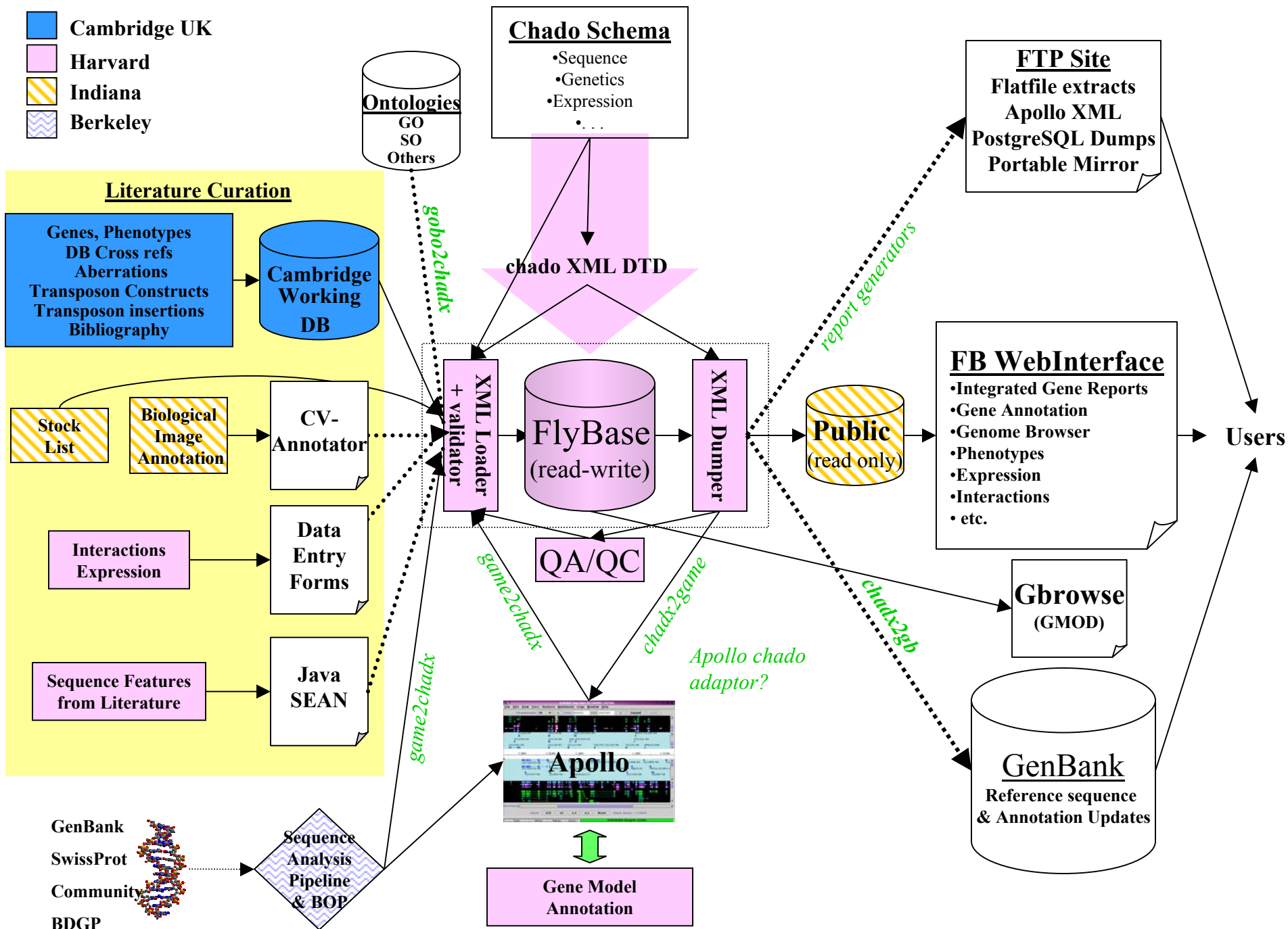
Query Performance

⌘ ROI Query

- ☑ Reasonable, not stellar performance on PostGreSQL with index on (srcfeature_id, min, max)
- ☑ Exploration of more sophisticated approaches yielded performance improvements in MySQL but not PostGreSQL
- ☑ PostGreSQL functions simplify queries, e.g. `select * from contains(src,min,max)`

⌘ Central Dogma Query

- ☑ a few seconds for 3 levels, all info
- ☑ 1 minute to include all overlapping features



XORT

XML Object to Relational Translator

⌘ Schema-driven tools

☒ DTD generator: DDL -> DTD

☒ also generates html, xml, .pl versions of schema

☒ Validator:

☒ Not connected:

- Syntax Verification: legal XML, correct element nesting
- Some Semantic verification: NULLness, cardinality, local ID reference

☒ Connected: reference validation

☒ Loader-only: constraints, triggers

☒ Loader: XML -> DB

☒ Dumper: DB -> XML

☒ driven by XML "dumpspec"

Mapping XML to R-DBMS

⌘ Policy#1: XML is independent of schema

- ☑ Pro: ensures modularity, freedom to change one without the other (but why would you want to?)
- ☑ Con: must maintain mapping when either changes / extends

⌘ Policy#2: XML locked to schema

- ☑ Pro: don't have to learn two things, mapping is frozen. Extension/revision is free.
- ☑ Con: see Pro above.

XORT Mapping



⌘ Elements

- ☒ Table
- ☒ Column (except primary key -- not visible in XML)

⌘ Attributes

- ☒ few and generic: transaction and reference control

⌘ Element nesting

- ☒ column within table
- ☒ joined table within table -- joining column is implicit
- ☒ foreign key table within foreign key column

⌘ Modules

- ☒ No module distinctions in chadoXML

⌘ Limitations of DTD

- ☒ Cardinality, NULLness, data type

```

<!-- ***** TABLE feature *****
*   feature_id      serial      not null      primary key
*   dbxref_id       int         foreign key to dbxref.dbxref_id - dbxref_id here is intended for
the primary dbxref for this feature.  Additional dbxref links are made via feature_dbxref
*   organism_id     int         not null      foreign key to organism.organism_id      unique(2)
*   name            varchar(255) - the human-readable common name for a feature, for display
*   uniquename      text        not null      unique(2) - the unique name for a feature; may not
be particularly human-readable
*   residues        text
*   seqlen          int
*   md5checksum     char(32)
*   type_id         int         not null      foreign key to cvterm.cvterm_id
*   timeaccessioned timestamp    not null      default current_timestamp -
timeaccessioned and timelastmodified are for handling object accession/ modification
timestamps (as opposed to db auditing info, handled elsewhere).  The expectation is that
these fields would be available to software interacting with chado.
*   timelastmodified timestamp    not null      default current_timestamp
***** -->
<!ELEMENT feature (dbxref_id | organism_id | name | uniquename | residues | seqlen |
md5checksum | type_id | timeaccessioned | timelastmodified)*
( feature_cvterm | analysis | featureloc | featureprop | analysisfeature |
feature_genotype | featurepos | feature_dbxref | feature_synonym | featurerange |
feature_relationship | feature_phenotype | feature_pub | interaction_subj |
interaction_obj | feature_expression)*
)>

<!ATTLIST feature
  id CDATA #IMPLIED
  ref CDATA #IMPLIED
  op (lookup | insert | update | force | delete) #IMPLIED>
<!ELEMENT dbxref_id ( #PCDATA | dbxref ) >      <!ATTLIST dbxref_id op (update) #IMPLIED>
<!ELEMENT organism_id ( #PCDATA | organism ) >  <!ATTLIST organism_id op (update)
#IMPLIED>
<!ELEMENT name #PCDATA >      <!ATTLIST name op (update) #IMPLIED>
<!ELEMENT uniquename #PCDATA >      <!ATTLIST uniquename op (update) #IMPLIED>
<!ELEMENT residues #PCDATA >      <!ATTLIST residues op (update) #IMPLIED>
<!ELEMENT seqlen #PCDATA >      <!ATTLIST seqlen op (update) #IMPLIED>
<!ELEMENT md5checksum #PCDATA >      <!ATTLIST md5checksum op (update) #IMPLIED>
<!ELEMENT type_id ( #PCDATA | cvterm ) >      <!ATTLIST type_id op (update) #IMPLIED>
<!ELEMENT timeaccessioned #PCDATA >      <!ATTLIST timeaccessioned op (update) #IMPLIED>
<!ELEMENT timelastmodified #PCDATA >      <!ATTLIST timelastmodified op (update) #IMPLIED>

```


Object References

**How to refer to persistent objects within XML?
(a.k.a. foreign key columns)**



⌘ By Unique Key Value(s)

- ☑ object can be in XML file or DB

⌘ By local ID

- ☑ only for references to objects in same XML file
- ☑ need not be in DB
- ☑ local ID can be any symbol - def before ref
- ☑ reduces duplication within XML

⌘ By Global accession

- ☑ currently only for feature
- ☑ simple extension mechanism using Perl fragments

Object Reference

by key values



```
<foreign_key_col>
  <primarytable>
    <keycol1>keyval1</keycol1>
    ... more key cols if needed
  </primarytable>
</foreign_key_col>
```

E.g.

```
<feature>
  <type_id>
    <cvterm>
      <cv_id>
        <cv>
          <name>Sequence Ontology</name>
        </cv>
      </cv_id>
      <name>exon</name>
    </cvterm>
  </type_id>
```

....

Object Reference by Local ID



```
<cv id="SO">  
  <name>Sequence Ontology</name>  
</cv>
```

```
<cvterm id="exon">  
  <cv_id>SO</cv_id>  
  <name>exon</name>  
</cvterm>
```

```
<feature>  
  <type_id>exon</type_id>
```

...

Object Reference

by Global Accession



```
<feature_relationship>  
  <subjfeature_id>GB:g012345  
  </subjfeature_id>
```

...

Transactions



⌘ **Lookup:** <table op="lookup">...

⌘ **Insert:** <table op="insert">...

⌘ **Delete:**

```
<table op="delete">
  <keycol1>val1</keycol1>...
```

⌘ **Update**

```
<table op="update">
  <keycol1>val1</keycol1>
  <keycol2>val2</keycol1>
  <keycol1>newval</keycol1>
```

⌘ **Force:** <table op="force">...

Combination of lookup, insert and update

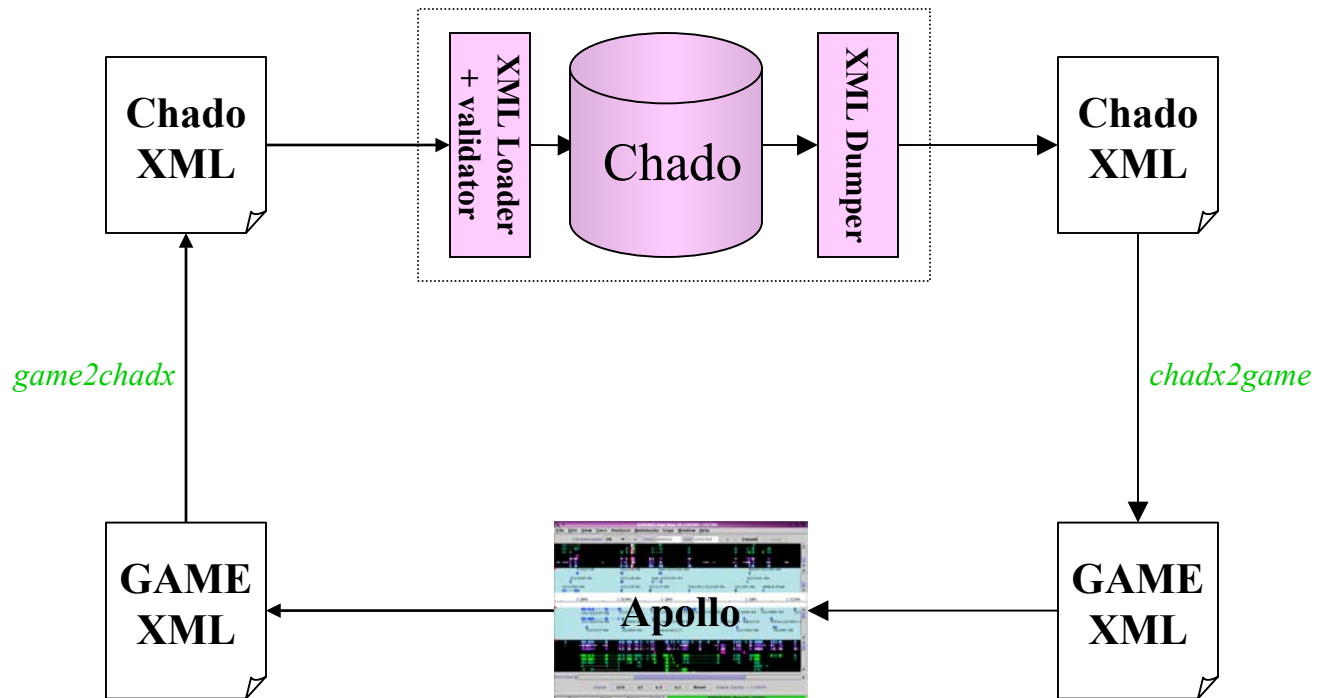
Dumper

XML-driven extraction



- ⌘ Default behavior: given an object class and ID, dump all direct values and linktables, with refs to foreign keys.
- ⌘ Nondefault behavior specified by XML dumpspecs using same DTD with a few additions:
 - ⊞ attribute dump= all | cols | select | none
 - ⊞ attribute test = yes | no
 - ⊞ element OR
 - ⊞ element _sql
 - ⊞ element _appdata
- ⌘ Workaround with views, _sql
- ⌘ Current use cases:
 - ⊞ Dump a gene for a gene detail page
 - ⊞ Dump a scaffold for Apollo

Chado <-> Apollo Interaction



DUMPER Concerns



- ⌘ Expressivity
- ⌘ Speed
- ⌘ XML file size
- ⌘ Memory

What's next



- ⌘ Debug Apollo / chado roundtrip
- ⌘ CV issues
 - ☑ Hierarchical queries
 - ☑ SO compliance
 - ☑ feature relationship types
- ⌘ Schema extensions
 - ☑ genetics module - review in Fall
 - ☑ expression
 - ☑ Function generally
- ⌘ UI development - not based on dumper

Architectural Principles



- ⌘ Semi-permeable XML layer:
 - ☑ support but do not require XML
- ⌘ Fix mapping, let schema vary
- ⌘ Plan for schema evolution -- schema-driven tools
- ⌘ Course-grained coupling of modules made possible by XML standardization

GMOD Architecture:

oxymoron?



- ⌘ Claim: self-contained, interoperable reusable complete system will not happen by itself
- ⌘ All tools speak all formats?
- ⌘ Standardize DB, API, or format?
 - ☑ Format standardization good for persistence, less good for querying
 - ☑ Bless existing format: ASN.1, GFF, BSML, SBML, ... or YAF?

Credits



Pinglei Zhou - loader, dumper, XML design

Frank Smutniak - game2chadx, chadx2game

Colin Wiel - gadfly2chado migration, schema

David Emmert - schema, migration

Chris Mungall, Suzi Lewis - schema, SO

Stan Letovsky - XML/tool design, dtd generator

Susan Russo, Mark Zytковicz - PostgreSQL

Don Gilbert - XML customer

Scott Cain - GBROWSE/Chado

Allen Day - schema (expression)

Hilmar Lapp - ROI query optimization

...