# And then magic happens…
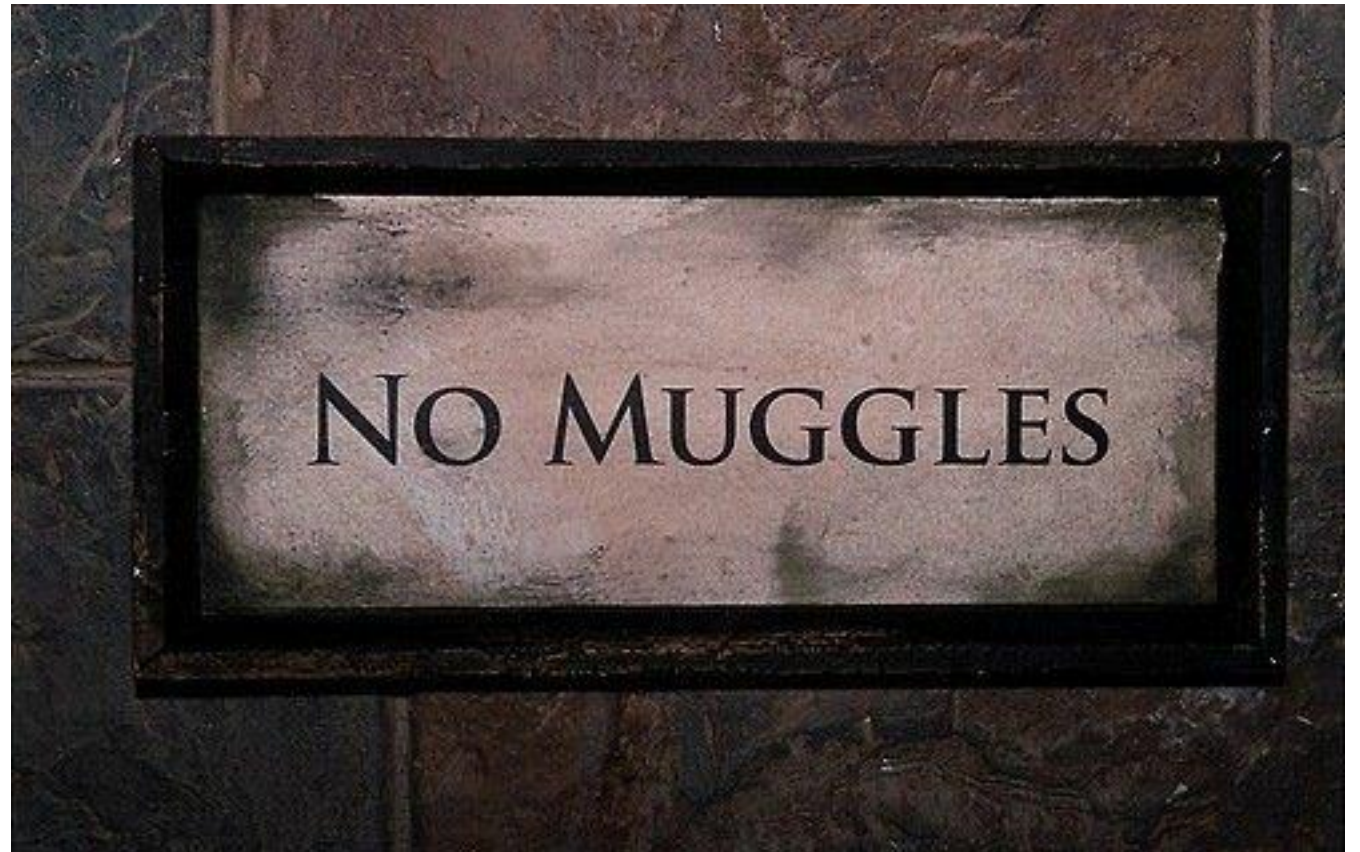
@Chris_Evelo

Maastricht University
WikiPathways team
ELIXIR Interoperability team
Open PHACTS

NO MUGGLES

# If we don't do the magic

# Lessons Learned from caBIG™

NCTR

January 21, 2010

Warren Kibbe

Northwestern University

wakibbe@northwestern.edu

# caBIG, the hard parts

o Progress is hard to measure

o Delivering exactly what people want is *not* the mission of caBIG

o Overspecifying, model-first, and semantic inflexibility are all enemies of caBIG's need to be agile and flexible

o The products are harder to use than anyone would like, but that is the nature of the space, so far…

# caBIG, revisited

o Fuzzy semantics and semantic web technologies will be embraced

o Complete specification is an unrealistic goal
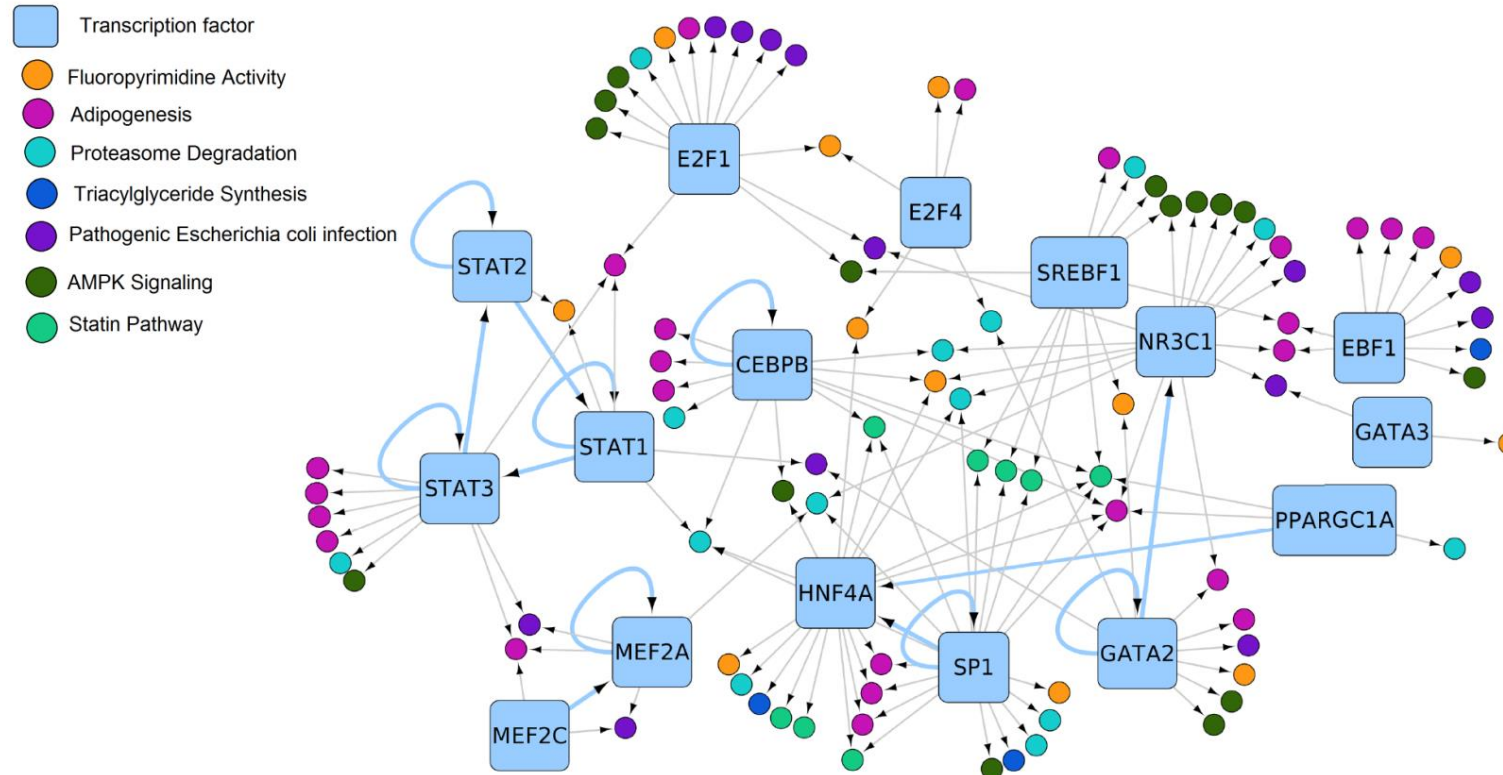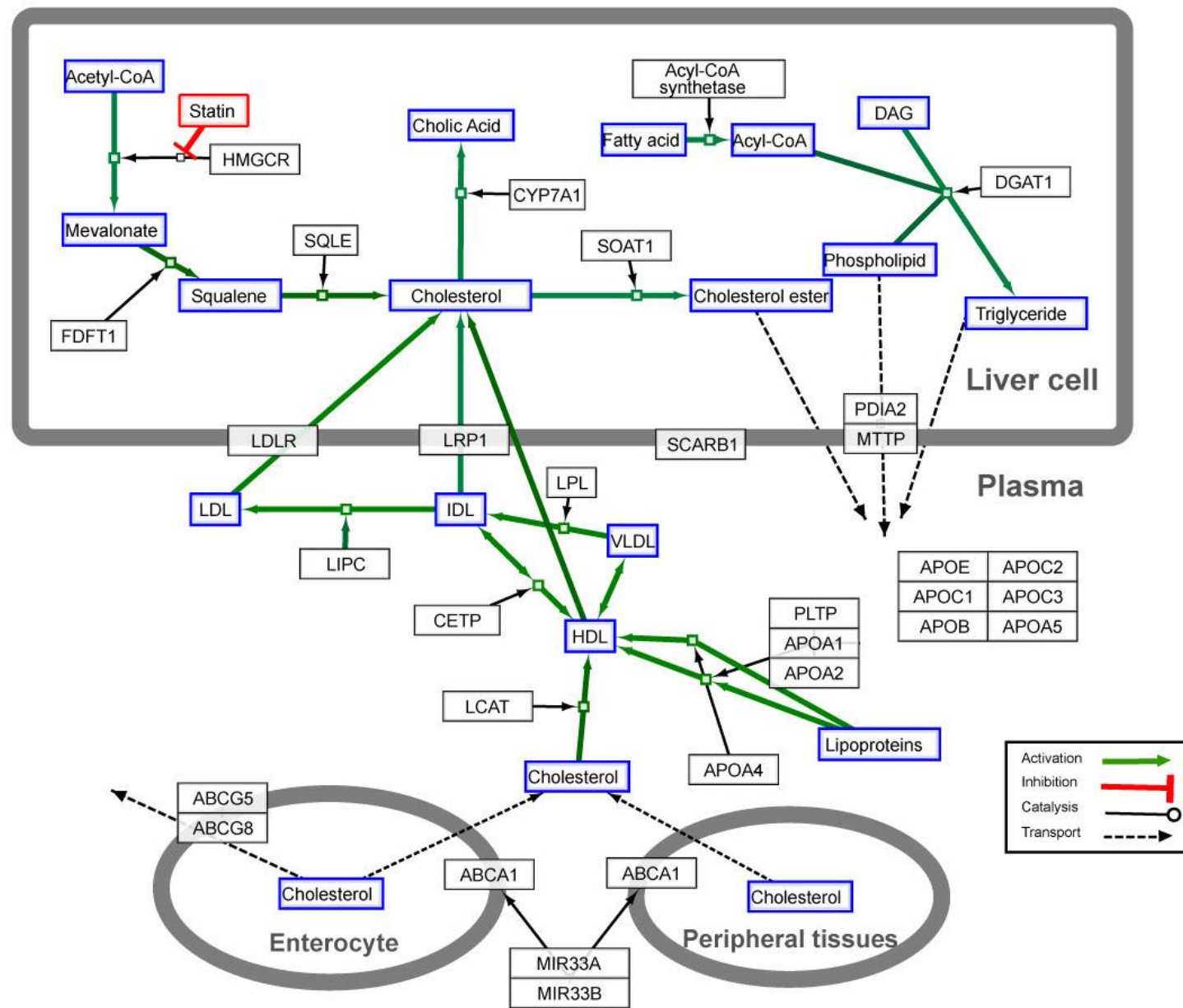
o Architect for change

Comfort Zone

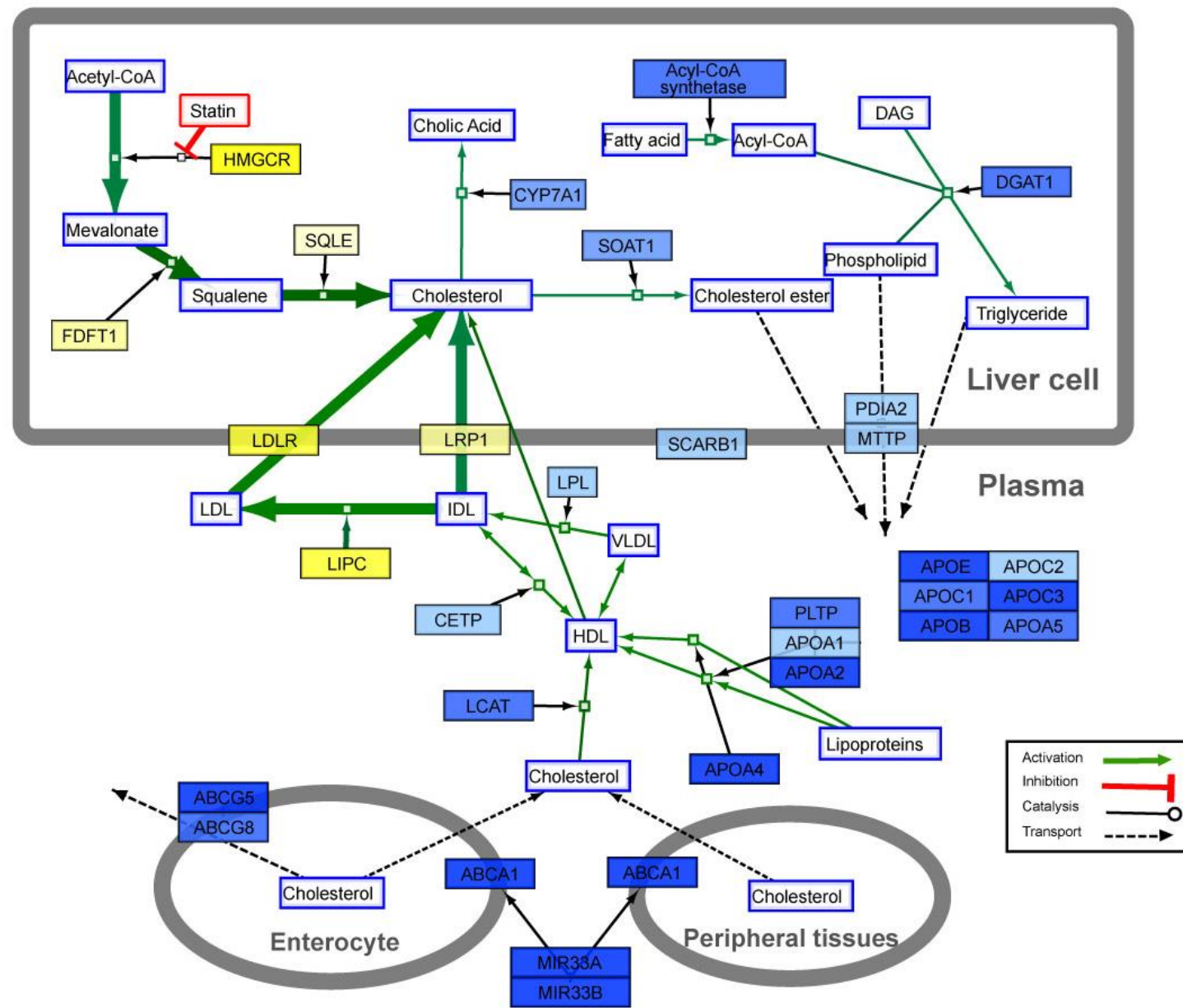Where the MAGIC Happens

2015 Jeannel King

# We can do things like this...
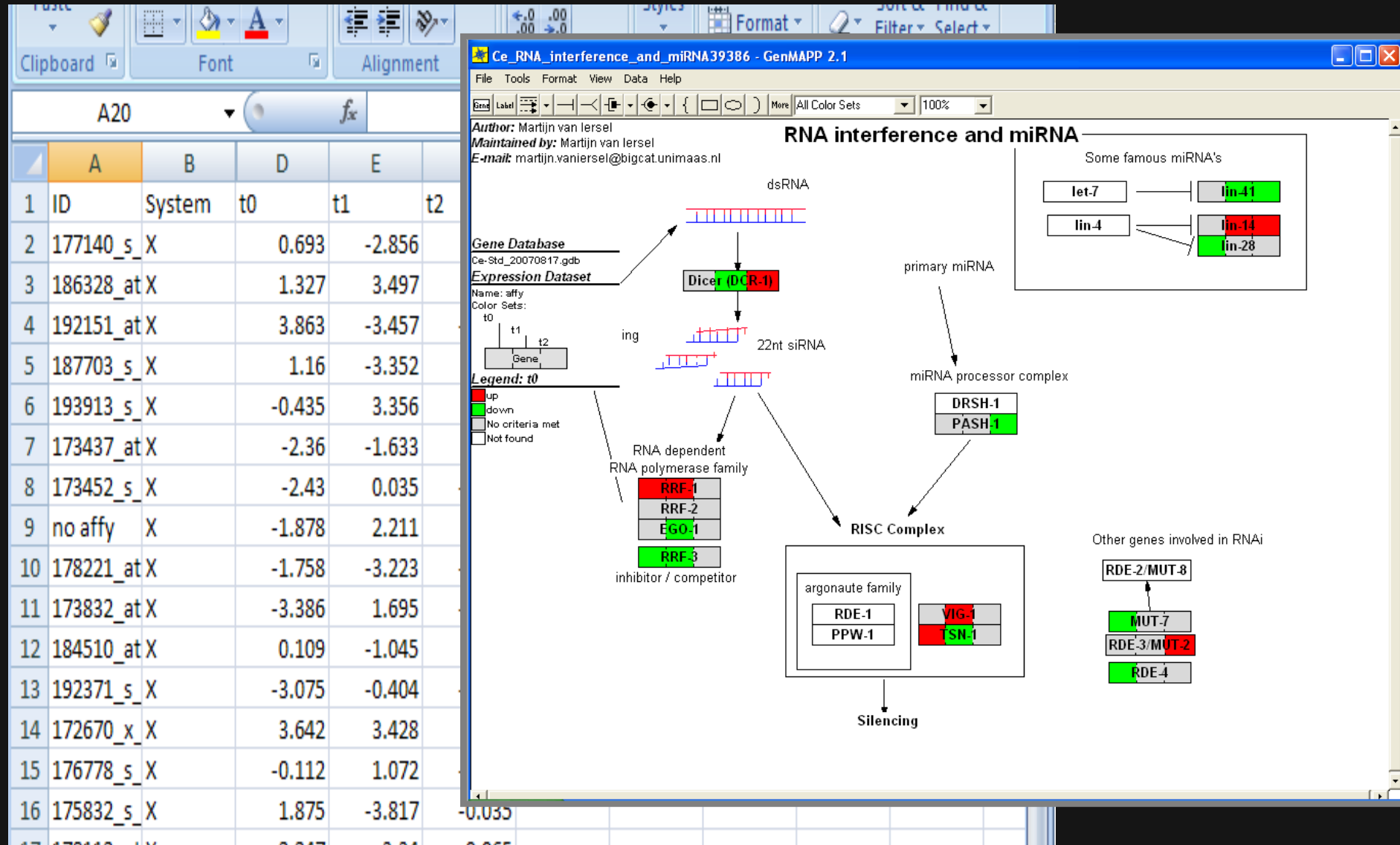


**Transcriptional regulation**

Using CyTargetLinker and the ENCODE proximal TF-target gene network, 16 transcription factors were identified in the selected pathways. Our analysis showed that TFs can be considered additional links between pathways and adding the regulatory interactions increases the overall connectivity of the network significantly.

# How to do data visualization?

# Connect to Genome Databases

# Connect to Genome Databases

# Backpages link to multiple databases

# You could do this for gene lists

# BridgeDb: Abstraction Layer



**class**
IDMapperRdb

*relational database*

**interface**
IDMapper

**class**
IDMapperFile

*tab-delimited text*

**class**
IDMapperBiomart

*web service*

# The magic is only in part in the tool

# And mostly in the content

# For database identifier mapping tools we have:

- A software framework (BridgeDb)
  - Application in WikiPathways, PathVisio, Cytoscape, R/Bioconductor
  - An installable webservice
  - Open source
  - Community based
  - Database based (small)
- A semantic web implementation (Open PHACTS IMS)
  - With installable Docker image
  - Linkset based (fast)
  - Does do transivity
- Identifiers.org for ID schema's and resolution

# For ID mappings we have:

- Gene product, ENSEMBL derived, databases for:
  - most health related species
  - Some bacteria
  - Some plants
- Metabolite database, HMDB & CheBi derived, for most "human" metabolites
- Reaction database from RHEA
- Linksets for all of these

# Other relevant:

- MyGene.info (Su lab)
  - also connected to WikiData
- DAVID (from Tigr)
  - recently updated (after years of inactivity)
  - often used (even in that period)
  - Should watch out for Goliath

# How do R&D companies use public data?



Patents  Literature  PubChem  Genbank

Databases

Downloads

Data Analysis

Data Integration

Firewalled Databases

# How do pharma companies use public data?

@gray_ala

# Apps

NO MUGGLES

| VoID | VoID | VoID | VoID | VoID |
|------|------|------|------|------|
| | Nanopub | | Nanopub | Nanopub |
| **RDF** | **RDF** | **RDF** | **RDF** | **RDF** |
| Db | Db | Db | Db | |

**Public Ontologies**

**Public Content** — **Commercial** — **User Annotations**

Linked Datasets as of August 2014

Publications
Life Sciences
Cross-Domain
Social Networking
Geographic
Government
Media
User-Generated Content
Linguistics

# Open PHACTS
## Open Pharmacological Space

# Apps

NO MUGGLES

| VoID | VoID | VoID | VoID | VoID |
| | Nanopub | | Nanopub | Nanopub |
| RDF | RDF | RDF | RDF | RDF |
| Db | Db | Db | Db | |

Public Ontologies

Public Content     Commercial     User Annotations

# Choose a standard

# Link one resource to another

# Or use both and map

# ID mapping in Open PHACTS:

- Semweb enabled (full URIs)
- Using linksets
- Transitivity (and limits for that)
  - gene -> protein -> has enzyme code
  - Protein -> has enzyme code -> other proteins

# This is not just Open PHACTS

Federated SPARQL queries:

e.g. find all genes related to disease, then all pathways with these genes…

Used as hackaton (swat4ls) examples/

Only works sometimes, by chance

Needs integrated ID mapping!

# Metabolite mapping needs

- More mappings! (plant products, drugs, xenobiotics)
- Ontology based mapping (CheBi)
- Because:
  - Palmitic acid is a fatty acid
  - R,R,R-tocopherol is a form of Vitamin E
- And these should (sometimes) map

# Also applies to biology:
# scientific lenses



"SoMeTiMes iWoNdeR wOuLdiT hAvE EnDeD diFFEreNtLy iF iChOsE AnOtHeR sToRy....?"

ZaRa
Sunday - Sep 23, 2012(1:30 pm)

# Concepts

- What we need is a big ARTA (also referred to as, synonyms) table
  - That allows typos..
- E.g. synonyms for things found by text mining
- Mapping to ontologies is only a partial solution

# ConceptWiki

Home   About

Search Search 🔍

![ConceptWiki logo] ConceptWiki

Search [_____] **Search 🔍**

"ConceptWiki is an editable community owned repository of concepts used to define all concepts unambiguously. **Anything you can write a Wikipedia article about is a concept.** The ConceptWiki terminology and identifiers can be freely downloaded and used to identify or refer to things in many different contexts."

2-oxoacetamide 5-Amino-2,3-dihydro-1,4-phthalazinedone ACETYLEUGENOL Anetholetrithione Chlor-Trimeton Allergy Contac 12 Hour Caplets Dehist Demilets Diathal (Veterinary) Dristan Sinus Eosine S 13 Midoxin Myocain Mytrate Omnidol Ornex Severe Cold Formula Para-phenacetin Paracetamolo Propadrine hydrochloride Restamin (TN) Ru-Tuss Tablets Syntodril TELDRIN dl-Chlorpheniramine maleate halquinol

© ConceptWiki 2010-2016

# It failed!

Because …

- Try it!
- Calling on a million minds for community annotation …
- But… people work on wikis because it gives them something right away

# And we still need it!

- Open PHACTS used an indexer
- These only know what is in there
- Doesn't for instance allow connecting text mining results

# Don't be afraid to reinvent wheels!

# Chemistry mapping

- Structure not ID based

- Allow substructure searches

- Open PHACTS open source ???

- We need it, may have to redo

# Ontology mapping

- Many available, even as services
- Often part of met resources
- Meta resources need feedback to original!

# Annotation tools use ontology terms

- But do not allow lookup
- Often (typically) do not show what ontology used
- Need frameworks

# Special cases

- Proteomics: peptides/fragments to proteins (Dutch proteomics meeting)

- Sequences: typical short sequence resolution (why BLAST the same things over and over?)  (Natasa Przulj at NetBio SIG)

# Mapping tools are core tools: need funding and sustainability

# Delivering

- Not what people want
- But what people need
- For what they want

- Where we will run into the next problem:

If we curate data well, and make all the technical connections real it may still be incomparable. Needs…. brains

People think that mathematics is complicated. Mathematics is the simple bit, it's the stuff we CAN understand. It's cats that are complicated.

— *John Horton Conway* —

Those who don't believe in magic will never find it.

Roald Dahl in "The Minpins"