# Clinical Genome Knowledge Base

# and Linked Data technologies

**Aleksandar Milosavljevic**

# Topics

1. ClinGen Resource project

2. Building the Clinical Genome Knowledge Base

3. Linked Data technologies

4. Using Linked Data technologies to enable pattern discovery in the Clinical Genome Knowledge Base

# Topics

1. ClinGen Resource project

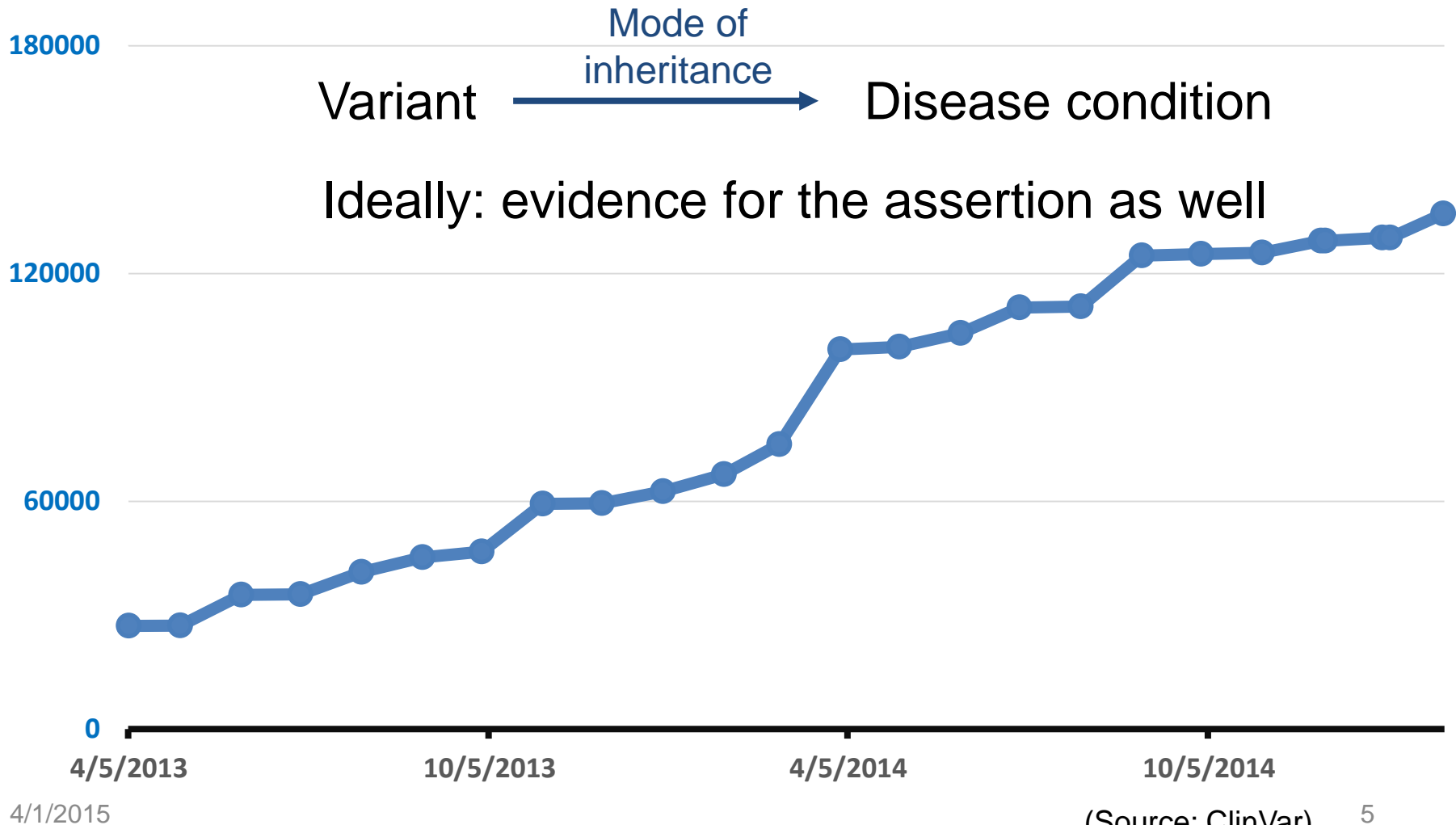1. Building the Clinical Genome Knowledge Base

2. Linked Data technologies

3. Using Linked Data technologies to enable pattern discovery in the Clinical Genome Knowledge Base
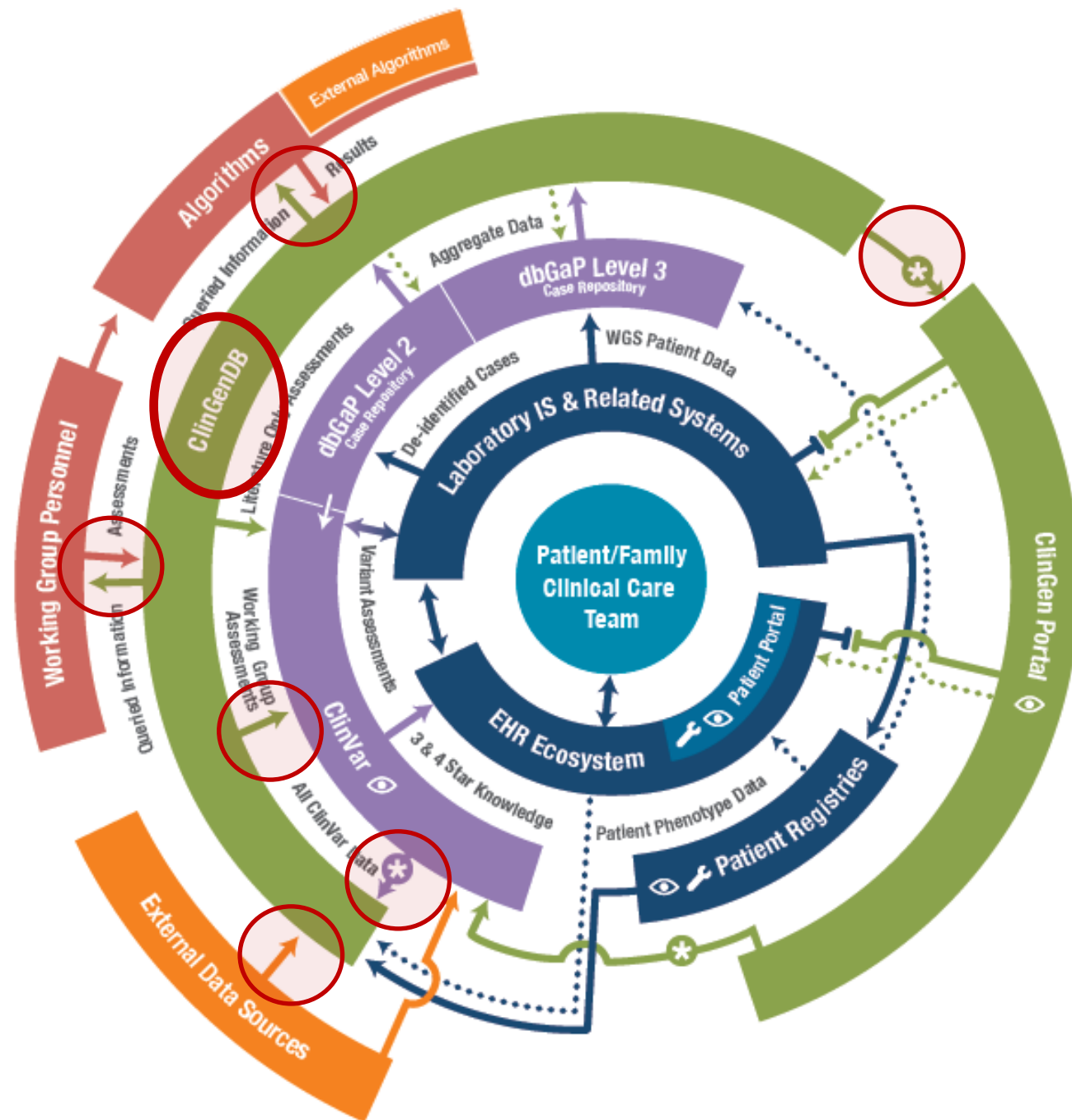
- Engage the clinical genomics community in data sharing efforts (sharing of data from genetic testing labs)

- Develop the infrastructure and standards to support curation of the knowledge about the role of genes and genetic variants in human diseases

- Incorporate machine-learning approaches to speed the identification of clinically relevant variants

# Variants in ClinVar
# (mostly from Diagnostic Genetic Testing)

Variant $\xrightarrow{\text{Mode of inheritance}}$ Disease condition

Ideally: evidence for the assertion as well

(Source: ClinVar)

# Topics

1. ClinGen Resource project

2. Building the Clinical Genome Knowledge Base

3. Linked Data technologies

4. Using Linked Data technologies to enable pattern discovery in the Clinical Genome Knowledge Base

# Evaluating pathogenicity of genetic variants for specific diseases

© American College of Medical Genetics and Genomics **ACMG STANDARDS AND GUIDELINES** | **Genetics in Medicine**

## Standards and guidelines for the interpretation of sequence variants: a joint consensus recommendation of the American College of Medical Genetics and Genomics and the Association for Molecular Pathology

Sue Richards, PhD[1], Nazneen Aziz, PhD[2,16], Sherri Bale, PhD[3], David Bick, MD[4], Soma Das, PhD[5], Julie Gastier-Foster, PhD[6,7,8], Wayne W. Grody, MD, PhD[9,10,11], Madhuri Hegde, PhD[12], Elaine Lyon, PhD[13], Elaine Spector, PhD[14], Karl Voelkerding, MD[13] and Heidi L. Rehm, PhD[15]; on behalf of the ACMG Laboratory Quality Assurance Committee

# Evaluating pathogenicity of genetic variants for specific diseases

## ACMG STANDARDS AND GUIDELINES

RICHARDS *et al* | Interpretation of sequence variants

**Table 3** Criteria for classifying pathogenic variants

| Evidence of pathogenicity | Category |
|---|---|
| Very strong | PVS1 null variant (nonsense, frameshift, canonical ±1 or 2 splice sites, initiation codon, single or multiexon deletion) in a gene where LOF is a known mechanism of disease |
| Strong | PS1 Same amino acid change as a previously established pathogenic variant regardless of nucleotide change<br>Example:    Val→Leu caused by either G>C or G>T in the same codon |
| Moderate | PM1 Located in a mutational hot spot and/or critical and well-established functional domain (e.g., active site of an enzyme) without benign variation |

| Evidence of benign impact | Category |
|---|---|
| Stand-alone | BA1 Allele frequency is >5% in Exome Sequencing Project, 1000 Genomes Project, or Exome Aggregation Consortium |
| Strong | BS1 Allele frequency is greater than expected for disorder (see **Table 6**) |

10

# ACMG STANDARDS AND GUIDELINES

| | Benign | | Pathogenic | | | |
|---|---|---|---|---|---|---|
| | **Strong** | **Supporting** | **Supporting** | **Moderate** | **Strong** | **Very strong** |
| **Population data** | MAF is too high for disorder BA1/BS1 OR observation in controls inconsistent with disease penetrance BS2 | | | Absent in population databases PM2 | Prevalence in affecteds statistically increased over controls PS4 | |
| **Computational and predictive data** | | Multiple lines of computational evidence suggest no impact on gene /gene product BP4<br><br>Missense in gene where only truncating cause disease BP1<br><br>Silent variant with non predicted splice impact BP7<br><br>In-frame indels in repeat w/out known function BP3 | Multiple lines of computational evidence support a deleterious effect on the gene /gene product PP3 | Novel missense change at an amino acid residue where a different pathogenic missense change has been seen before PM5<br><br>Protein length changing variant PM4 | Same amino acid change as an established pathogenic variant PS1 | Predicted null variant in a gene where LOF is a known mechanism of disease PVS1 |
| **Functional data** | Well-established functional studies show no deleterious effect BS3 | | Missense in gene with low rate of benign missense variants and path. missenses common PP2 | Mutational hot spot or well-studied functional domain without benign variation PM1 | Well-established functional studies show a deleterious effect PS3 | |
| **Segregation data** | Nonsegregation with disease BS4 | | Cosegregation with disease in multiple affected family members PP1 | Increased segregation data →→→ | | |
| **De novo data** | | | | De novo (without paternity & maternity confirmed) PM6 | De novo (paternity and maternity confirmed) PS2 | |
| **Allelic data** | | Observed in trans with a dominant variant BP2<br><br>Observed in cis with a pathogenic variant BP2 | | For recessive disorders, detected in trans with a pathogenic variant PM3 | | |
| **Other database** | | Reputable source w/out shared data = benign BP6 | Reputable source = pathogenic PP5 | | | |
| **Other data** | | Found in case with an alternate cause BP5 | Patient's phenotype or FH highly specific for gene PP4 | | | |

11

# Evaluating pathogenicity of genetic variants for specific diseases

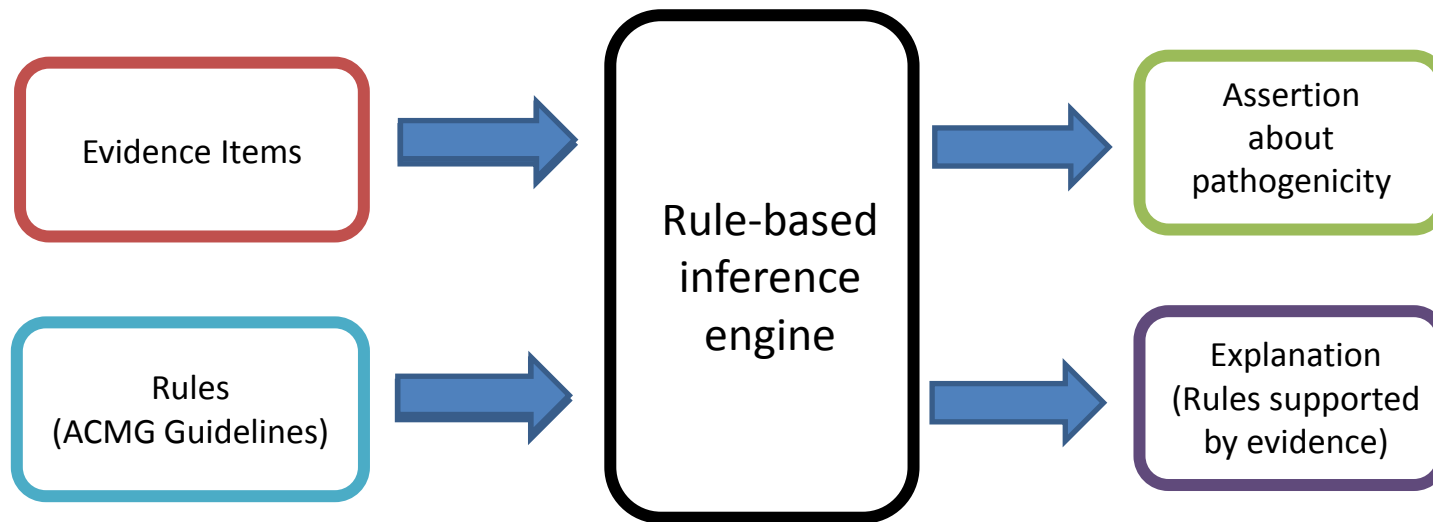**Table 5** Rules for combining criteria to classify sequence variants

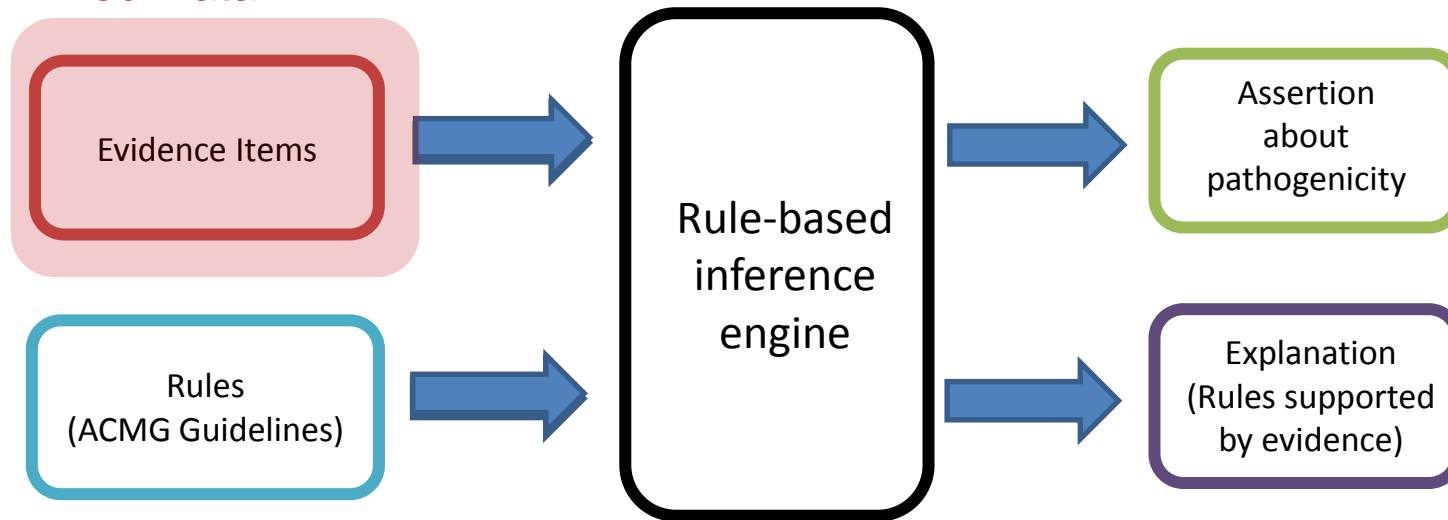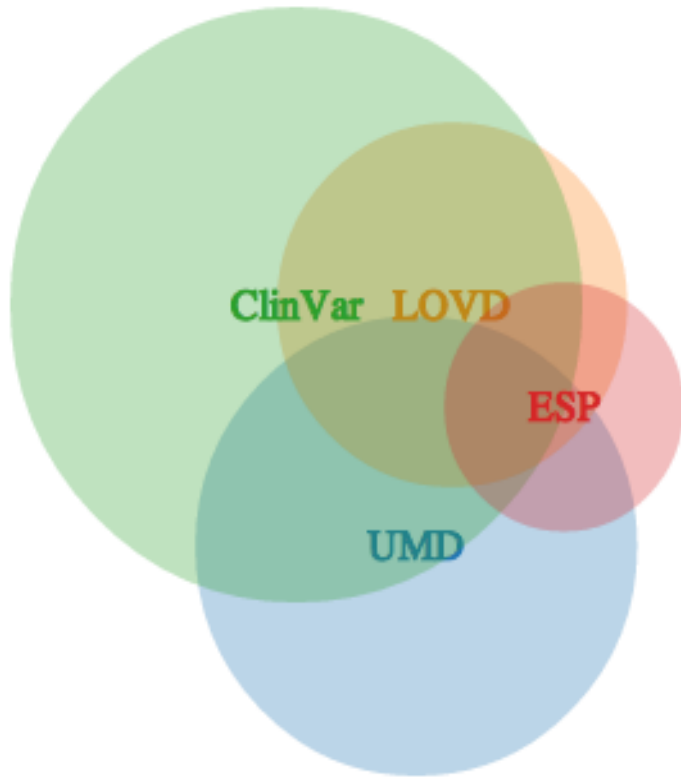| Pathogenic | (i) 1 Very strong (PVS1) *AND* |
|---|---|
| | (a) ≥1 Strong (PS1–PS4) *OR* |
| | (b) ≥2 Moderate (PM1–PM6) *OR* |
| | (c) 1 Moderate (PM1–PM6) and 1 supporting (PP1–PP5) *OR* |
| | (d) ≥2 Supporting (PP1–PP5) |
| | (ii) ≥2 Strong (PS1–PS4) *OR* |

| Benign | (i) 1 Stand-alone (BA1) *OR* |
|---|---|
| | (ii) ≥2 Strong (BS1–BS4) |

# Evaluating pathogenicity of genetic variants for specific diseases

Evidence Items → Rule-based inference engine → Assertion about pathogenicity

Rules (ACMG Guidelines) → Rule-based inference engine → Explanation (Rules supported by evidence)

# Evaluating pathogenicity of genetic variants for specific diseases

# Evidence code PM1: mutational hotspot

| Moderate | PM1 Located in a mutational hot spot and/or critical and well-established functional domain (e.g., active site of an enzyme) without benign variation |
|---|---|

**Variant frequency profile in normal population**



**Cancer-predisposing variants**

# Collating information about variants from public sources



**BRCA1**

ClinVar  LOVD

ESP

UMD

**APC**

dbSNP  ClinVar

LOVD

UMD

# Ontology traversal to collate variants linked to related disease conditions



**Inferred higher-level terms**

**Original disease data**

**Gene:TP53**

Reproductive organ cancer

Integumentary system cancer

Gastrointestinal system cancer

Respiratory system cancer

# Topics

1. ClinGen Resource project

2. Building the Clinical Genome Knowledge Base

3. Linked Data technologies

4. Using Linked Data technologies to enable pattern discovery in the Clinical Genome Knowledge Base

# RDF: Resource Description Framework

- Originally: metadata data model describing web resources

- Web resource = anything identifiable by a URL

- Subject-predicate-object expressions
  (related to classical entity-property-value)

- Subject = anything identifiable by a URL / URI

- Expression = edge in a "Knowledge Graph"

- Wkipathways "serializes" biological pathway data as RDF

Example:
Metastatic brain tumor pathway



19

# RDFs: RDF Schema

Example:
classification of body fluids

**Uses RDF to express data model (data schema)**

**Some RDFs predicates**:

- subClassOf – the subject is a subclass of a class

- subPropertyOf – the subject is a subproperty of a property

- domain – domain of the subject property

- range - range of the subject property

# OWL: Web Ontology Language

Example: "ontology traversal"
to annotate variants
to more general disease categories

- Uses RDF to express to describe taxonomies and classification networks

- Defines the structure of knowledge for various domains

# Some Linked Data projects

**"RDF-like graphs":**

- Facebook Social Graph -- Open Graph API
- Google Knowledge Graph / Vault
- WikiData

**RDF graphs:**

- data.gov
- Bio2RDF
- Bioontologies / Bioportal -- RDF, RDFs, OWL

# What is Linked *Open* Data?

# Linked RDF Standards

1999-2014      RDF, RDFs, OWL  *-- graph syntax and semantics*

2014           JSON-LD *-- JSON -- data format for RDF graphs*

Feb 26 2015   Linked Data Platform 1.0 – *HTTP operations*

- Use URIs as names for things

- Use HTTP URIs so that people can look up those names

- When someone looks up a URI, provide useful information, using the standards (RDF*, SPARQL)

- Include links to other URIs, so that they can discover more things

# Topics

1. ClinGen Resource project

2. Building the Clinical Genome Knowledge Base

3. Linked Data technologies

4. Using Linked Data technologies to enable pattern discovery in the Clinical Genome Knowledge Base

# Use Case #1: Discovering "match" patterns

# Use Case #2: Discovering Recurrent Mutually Exclusive (RME) mutational patterns in cancer

# Hallmark phenotypes of cancer are acquired by positive selection

(Hanahan and Weinberg, Cell 2001 and 2011)



Figure 1. Acquired Capabilities of Cancer

# Recurrent mutually exclusive mutational patterns identified in key glioblastoma pathways

**Comprehensive genomic characterization defines human glioblastoma genes and core pathways** TCGA Research Network, *Nature* **455**, 1061-1068 (2008)
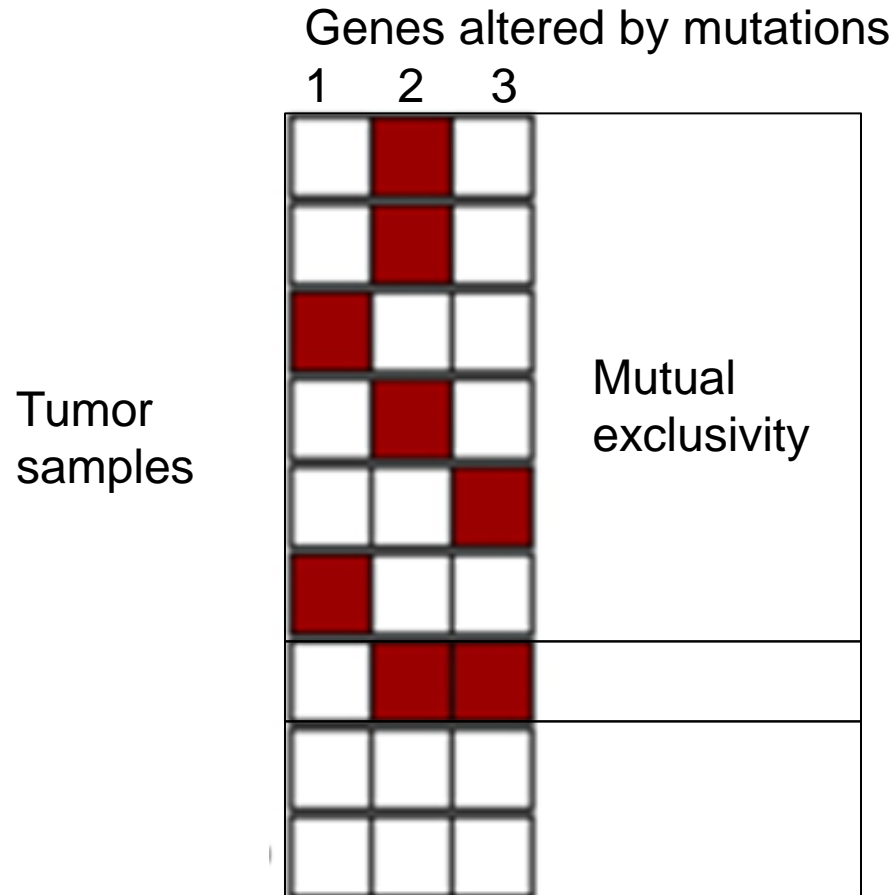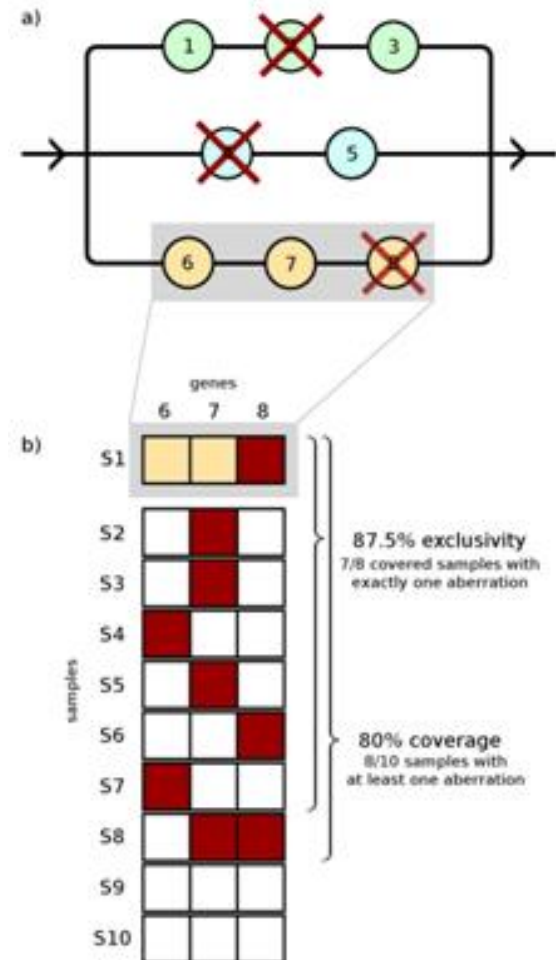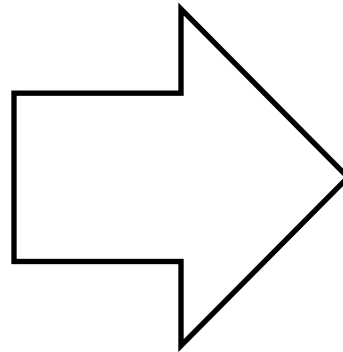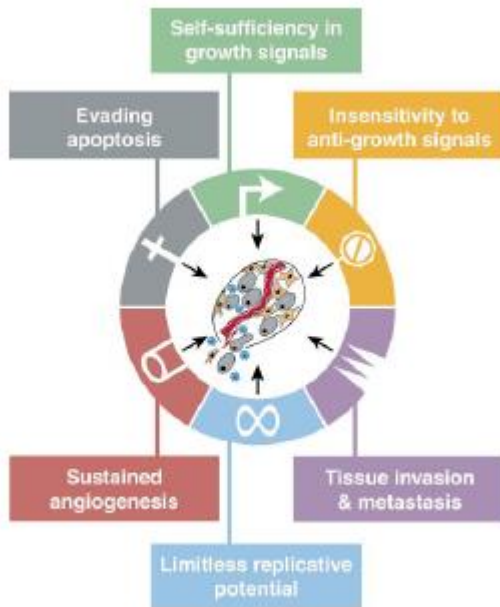
# Aristotle's Zoology: either horns or tusks

…[Aristotle] believed that **nature was economical** and gave no animal too many gifts, observing that **no animal possessed both horns and tusks**

| Acquired capability of animals: defense | |
|---|---|
| **Horns** | **Tusks** |
| |  |
|  | |
| |  |

# Discovering recurrent mutually exclusive patterns

Genes altered by mutations

Tumor samples

Mutual exclusivity

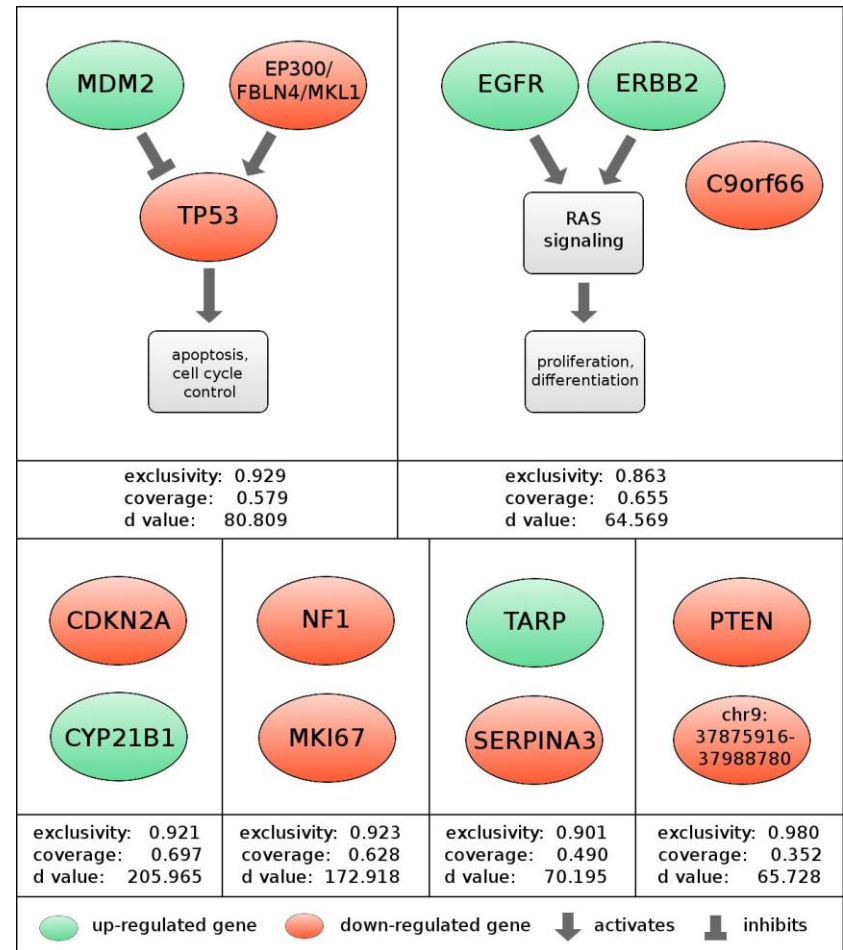# Discovering recurrent mutually exclusive patterns



(Hanahan and Weinberg, Cell 2001 and 2011)

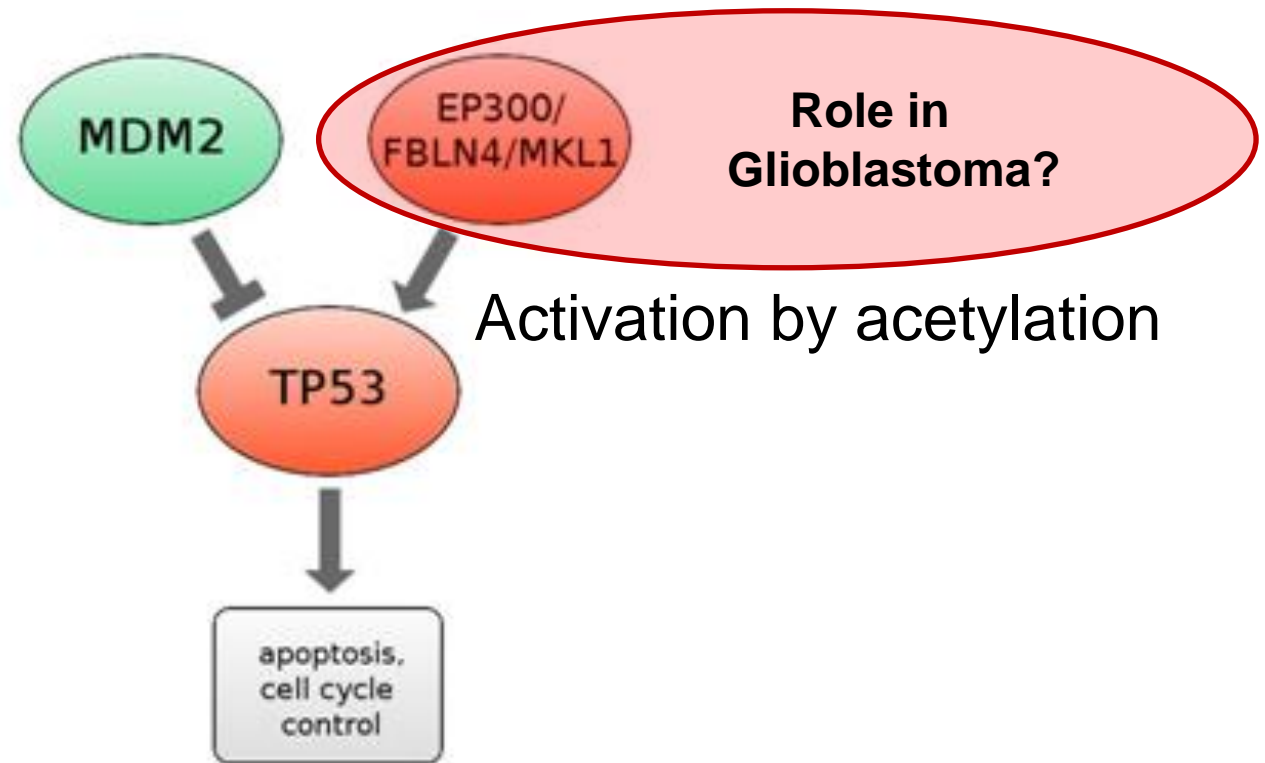# Can key modules / pathways be discovered based on RME patterns alone ?

(Miller CA et al. BMC Med Genomics. 4(1):34 2011)

# RME algorithm applied to glioblastoma TCGA collection (145 samples)

Rediscovered modules within all 3 pathways in glioblastoma reported by the TCGA Research Network, *Nature* **455**, 1061-1068 (2008).
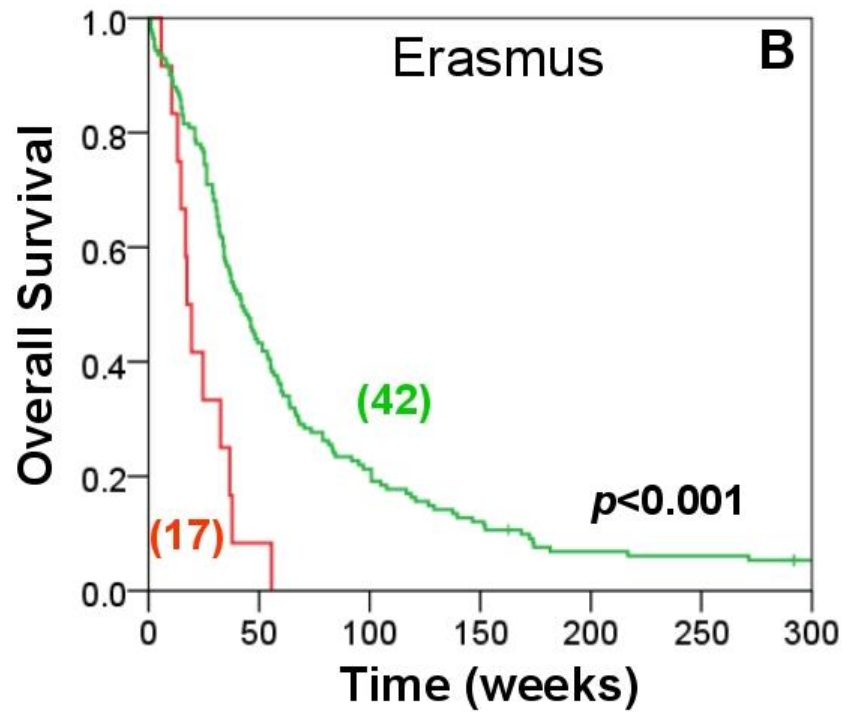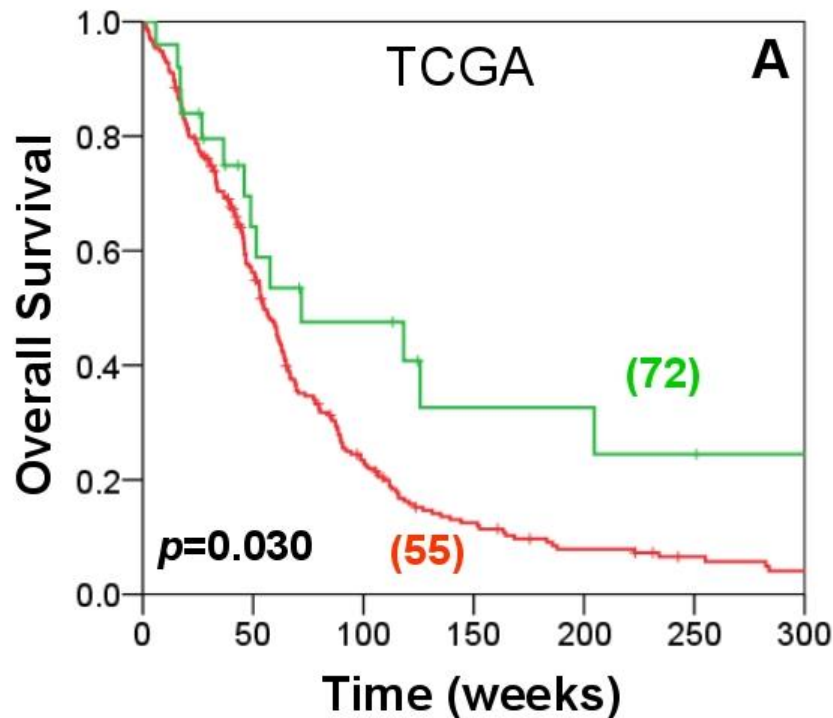
(Miller CA et al. BMC Med Genomics. 4(1):34 2011)

# RME algorithm applied to glioblastoma TCGA collection (145 samples)



exclusivity: 0.929
coverage: 0.579
d value: 80.809

# Effect of *EP300* expression
# on grade-independent and age-independent
# survival in glioblastoma

(Miller CA et al. BMC Med Genomics. 4(1):34 2011)

# Use Case #2: Discovering Recurrent Mutually Exclusive (RME) mutational patterns in cancer

Step 1: Identify variants associated with a specific phenotype (cancer)
> Disease ontology traversal

Step 2: Find groups of genes that form RME pattern

Step 3: Examine connectedness of the genes within pathways/networks
> Evaluate pattern of connectedness within networks or pathways

# Topics

1. ClinGen Resource project

2. Building the Clinical Genome Knowledge Base

3. Linked Data technologies

4. Using Linked Data technologies to enable pattern discovery in the Clinical Genome Knowledge Base

# Conclusion

The **Web of hyperlinks** (Web 1.0, 2.0) has greatly amplified the impact of the Human Genome Project
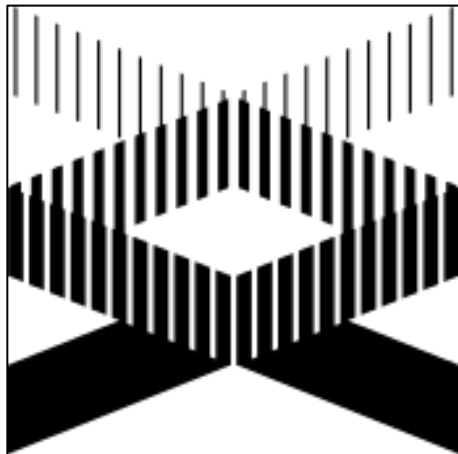
The **Web of Linked Data** (Web 3.0) will enable creation of a computable Clinical Genome Knowledge Base to inform

- Clinical interpretation of genomic variation

- Pattern discovery leading to new hypotheses

# Acknowledgements

**Baylor College of Medicine**

Sharon Plon
Andrew R. Jackson
Xin Feng
Ronak Patel
Rajarshi Ghosh

ClinGen
Clinical Genome Resource

National Human Genome Research Institute