**2$^e$ édition des Journées Bioss-IA**

# Search of Therapeutic Targets on the Hepatocellular Carcinoma with Database Extraction and Graph Coloring Methods

**Recherche de cibles thérapeutiques pour le carcinome hépatocellulaire à l'aide d'extraction de bases de données et de méthodes de coloration de graphes**

## Maxime FOLSCHETTE

**Current occupation:**
CNRS & Institut Français de Bioinformatique (IFB)
Laboratoire des Sciences du Numérique de Nantes (LS2N)
groups: GDD & COMBI

**Previous occupation:**
Université de Rennes 1
IRISA & IRSET
groups: Dyliss & Dymec

maxime.folschette@ls2n.fr
http://maxime.folschette.name/

2018-12-19

Context

**Hepatocellular carcinoma** **(HCC)**

- Most widespread liver cancer, $3^{rd}$ most deadly cancer
- Mainly provoked by hepatitis and fibrosis
- Late diagnosis and difficult to treat (resection, transplant, ablation, chemo-embolization)
- Very low survival rate: from weeks to months

**Objective**

- Gather data from **ICGC**
  [Hudson and The International Cancer Genome Consortium, 2010]
- Distinguish two tumor stages: early stage vs. late (invasive) stage
- Watch expression change (differential expression)

https://dcc.icgc.org/projects/LIHC-US

# LIHC-US

---

**LIHC-US in ICGC**

Project for liver HCC (USA)

- 294 samples with gene expression data
- Primary tumor on solid tissue only
- 20502 genes
- 16282 genes when excluding low expression

But no tumor grade annotation!

⇒ We need a **criterion** to distinguish tumor stages

---

**Objectives**

1) Clustering on the **criterion** ⇒ Two groups
2) Differential analysis on the two groups

---

## Epithelial-Mesenchymal Transition



**Epithelial-mesenchymal transition** **(EMT)**

**Epithelial cells**
Adhesive

**Mesenchymal cells**
Motile & invasive

EMT

- De-differentiation of epithelial cells to mesenchymal cells
- Gain ability to **remodel** the extra-cellular matrix and **migrate**
- Invasive cancer cells ⇒ **metastasis**

# Epithelial-Mesenchymal Transition

**Epithelial-mesenchymal transition** **(EMT)**



- De-differentiation of epithelial cells to mesenchymal cells
- Gain ability to **remodel** the extra-cellular matrix and **migrate**
- Invasive cancer cells ⇒ **metastasis**
- Indication of **tumor stage** ⇒ **Criterion**

# Epithelial-Mesenchymal Transition



**Epithelial-mesenchymal transition** **(EMT)**

**Epithelial cells**
Adhesive

**Mesenchymal cells**
Motile & invasive

EMT →

**Early stage**　　　　**Late stage**

- De-differentiation of epithelial cells to mesenchymal cells
- Gain ability to **remodel** the extra-cellular matrix and **migrate**
- Invasive cancer cells ⇒ **metastasis**
- Indication of **tumor stage** ⇒ **Criterion**
- **EMT signature** = Set of genes that are over-expressed during EMT (includes TGF-β)
- Downloaded on **GSEA** [Subramanian et al., 2005]

http://software.broadinstitute.org/gsea/msigdb/cards/HALLMARK_
EPITHELIAL_MESENCHYMAL_TRANSITION.html

http://software.broadinstitute.org/gsea/msigdb/cards/HALLMARK_
EPITHELIAL_MESENCHYMAL_TRANSITION.html

## Workflow of the Project

ICGC expression data
Clustering on EMT signature

↓

Differential expression analysis
→ Interesting genes

↓

Extraction of the pathways from
Kegg (Stream)

↓

Spread coloring and
make predictions (Iggy)

↓

Robustness analysis

294 samples (LIHC-US)

Group $A$ = Low expression of the EMT signature
Group $C$ = High expression of the EMT signature

## Workflow of the Project



**ICGC** expression data
Clustering on EMT signature

2 groups

Differential expression analysis
$\rightarrow$ Interesting genes

Extraction of the pathways from
**Kegg** (**Stream**)

Spread coloring and
make predictions (**Iggy**)

Robustness analysis

Differential Analysis

**Fold-change definition**

- Consider groups $A$ (lowest expression of EMT) and $C$ (resp. highest)
- For each gene $g$, compute mean value for group $A$ (resp. $C$)
- Differential analysis:

$$\text{fold-change}(g) = \text{mean}_g(C) \ / \ \text{mean}_g(A)$$

## Selected genes

**Criteria**

- Adjusted P-value $< 10^{-5}$
- $\log_2(\text{fold-change}) > 2$       (up-regulated genes)
- $\log_2(\text{fold-change}) < -0,5$    (down-regulated genes)

**Selected genes**

- 821 up-regulated genes
- 1092 down-regulated genes
  - $= 1913$ genes

**Objectives**

1) Extract a graph from **Kegg** [Kanehisa et al., 2017] using these genes, with the tool **Stream**
2) Coloring + predictions with **Iggy** [Thiele et al., 2015]

## Workflow of the Project



**ICGC** expression data
Clustering on EMT signature

2 groups

↓

Differential expression analysis
→ Interesting genes

$\simeq$ 2'000 genes

↓

Extraction of the pathways from
**Kegg** (**Stream**)

↓

Spread coloring and
make predictions (**Iggy**)

↓

Robustness analysis

## Pathway Commons + Bravo

### Pathway Commons [Cerami et al., 2010]

- A gathering of **25 pathway databases**
- Contains: PID, Kegg, Reactome, CTD, Panther, …
- Common ontology (BioPAX)
- Freely available
- SPARQL endpoint

### BRAvo [Lefebvre et al., 2017]

- **Interrogates** Pathway Commons with SPARQL queries
- Search and fusion of synonyms, optimizations
- (Incomplete) visualization tool

## Problems with Pathway Commons

**Problems with Pathway Commons**

- **Very heterogeneous data**, curation depends on the data sources
- The BioPAX ontology is big and difficult to use
- Unification must be done by the user, based on gene unames (fast) or identifies (slow)
- Updated without history (bad for reproducibility)

**Problems with BRAvo**

- Still in development, only regulation at the time (no signaling)
- Struggles with the heterogeneous content of Pathway Commons
- Unification was still incomplete

# Kegg + Stream

## Kegg [Kanehisa et al., 2017]

- Homogeneous data

- Categories: 2. Genetic Information Processing
    3. Environmental Information Processing
    4. Cellular Processes
    5. Organismal Systems

- Already formatted and curated by Arnaud Poret

$$A \xrightarrow{+/-} B$$

SIF format:　　　　　　　　　　　　　　　　"$A$ positively/negatively influences $B$"

  - Genes (XXX_gen)
  - Proteins (XXX_prot)
  - Complexes (XXX::YYY::ZZZ)

## Stream (Arnaud Poret)

- Ad-hoc program for upstream graph extraction

- Extract the part of the graph for which we have expression data (25%)

**Graph content:**

- 3'383 nodes
- 13'771 edges
  - — 11'661 activations
  - — 2'110 inhibitions

1913 genes from the differential expression
**Only 209 are found in Kegg:**

- ■ 138 up-regulated
- ■ 71 down-regulated
- ■ 3174 new nodes

Nodes with up to:
**92 incoming influences**
**79 outgoing influences**
$\rightarrow$ Nodes with a lot of impact on the network

## Workflow of the Project



**ICGC** expression data
Clustering on EMT signature

2 groups

Differential expression analysis
→ Interesting genes

$\simeq$ 2'000 genes

Extraction of the pathways from
**Kegg** (**Stream**)

$\simeq$ 3'400 nodes, 14'000 edges

Spread coloring and
make predictions (**Iggy**)

Robustness analysis

# Graph Coloring

- Coloring = information attached to nodes about over- or under-expression
  
  $(X)$ = over-expressed    $(Y)$ = under-expressed

- Provenance = experimental (expression data) & computational (inference)



Given by the
experimental data

# Graph Coloring

- Coloring = information attached to nodes about over- or under-expression

  $X$ = over-expressed     $Y$ = under-expressed

- Provenance = experimental (expression data) & computational (inference)
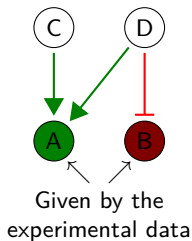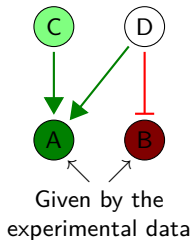


Given by the
experimental data

# Graph Coloring

- Coloring = information attached to nodes about over- or under-expression

  $X$ = over-expressed       $Y$ = under-expressed

- Provenance = experimental (expression data) & computational (inference)



**Consistent**

# Graph Coloring

- Coloring = information attached to nodes about over- or under-expression

  $\left(X\right)$ = over-expressed    $\left(Y\right)$ = under-expressed

- Provenance = experimental (expression data) & computational (inference)



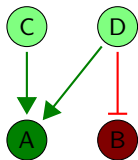**Consistent**          **Consistent**
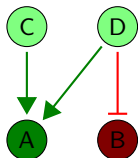
# Graph Coloring

- Coloring = information attached to nodes about over- or under-expression

  (X) = over-expressed     (Y) = under-expressed

- Provenance = experimental (expression data) & computational (inference)



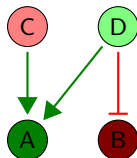**Consistent**          **Consistent**          **Inconsistent**

# Graph Coloring

- Coloring = information attached to nodes about over- or under-expression

  (X) = over-expressed    (Y) = under-expressed

- Provenance = experimental (expression data) & computational (inference)



**Consistent**      **Consistent**      **Inconsistent**      **Inconsistent**

# Graph Coloring

- Coloring = information attached to nodes about over- or under-expression

  (X) = over-expressed      (Y) = under-expressed

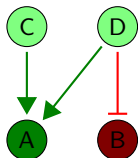- Provenance = experimental (expression data) & computational (inference)
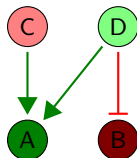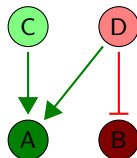


**Consistent**         **Consistent**         Inconsistent         Inconsistent
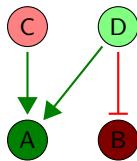
- Compute all colorings without inconsistencies
- **Prediction** = a node that is always colored the same

  Here, only 1 prediction: (D)

- All computed by **Iggy** [Thiele et al., 2015] (Answer Set Programming)

# Trivial Predictions

---

**"Trivial" prediction**

- Protein predicted the same as its observed gene
- Rarely brings new information
- Useful for validation

---



$\leftarrow$ Observations (inputs) $\rightarrow$

$\leftarrow$ Predictions (result) $\rightarrow$

**Knowledge from experiments:**

- 138 up-regulated
- 71 down-regulated

**Computational predictions:**

- 92 predicted $(+)$
  - 24 non-trivial
- 54 predicted $(-)$
  - 33 non-trivial

**70% more information** compared to only knowledge from experiments

# Computational predictions (results of Iggy)

# Workflow of the Project



**ICGC** expression data
Clustering on EMT signature

2 groups

Differential expression analysis
→ Interesting genes

$\simeq$ 2'000 genes

Extraction of the pathways from
**Kegg** (**Stream**)

$\simeq$ 3'400 nodes, 14'000 edges

Spread coloring and
make predictions (**Iggy**)

Predictions:
92 (+) and 54 (−)

Robustness analysis

- 209 inputs

**Matching between comp[al] predictions and ICGC expression data:**

- 124 match
  - 36 non-trivial
- 17 do **not** match
  - 16 non-trivial
- 5 not found in ICGC data

**88% matching**

**69% non-trivial**

$\rightarrow$ Good overlap

# Cross-Validation

## Sampling

- Consider a range of samplings (10%, 15%, 20%, ... 95%)
- Randomly pick x% of under- and over-expressed genes (observations)
- Compute the predictions on this sample ; repeat 100 times

### Score compared to the original data

- Compare the predictions to the original ICGC data
- Give a score to each set of predictions
  - → Scores converge to the final score at 100%

### Robustness of the prediction of each node

- Compare the predictions to the final sampling of 100%
  - → Not a lot of variability in the prediction types → Robust

Boxplot of the scores for each sampling & curves of the number of predictions

## Prediction Results

---

**New results compared to ICGC : complexes**

**Complexes predicted:**

- NFKB1::BCL3 ($+$)

- NFKB2::RELB ($+$)

- JUND::NACA ($-$)

---

**Results conflicting with ICGC data**

**Computational predictions which are different from differential analysis:**

- BAK1_gen, BMP4_gen, CREB1_prot, EIF4EBP2_prot, IGFBP3_gen,
  IGFBP3_prot, NR0B2_gen, NR0B2_prot, NR1H4_gen, NR1H4_prot,
  NR3C2_gen, NR3C2_prot, SESN3_gen, SESN3_prot, THBS1_gen,
  TNFRSF10A_gen, TP53_prot

---

## Prediction Results

---

**New results compared to ICGC : complexes**

**Complexes predicted:**

- NFKB1::BCL3 ( + )

- NFKB2::RELB ( + )

- JUND::NACA ( − )

---

**Results conflicting with ICGC data**

**Computational predictions which are different from differential analysis:**

- BAK1_gen, BMP4_gen, CREB1_prot, EIF4EBP2_prot, IGFBP3_gen, IGFBP3_prot, NR0B2_gen, NR0B2_prot, NR1H4_gen, NR1H4_prot, NR3C2_gen, NR3C2_prot, SESN3_gen, SESN3_prot, THBS1_gen, TNFRSF10A_gen, **TP53_prot**

---

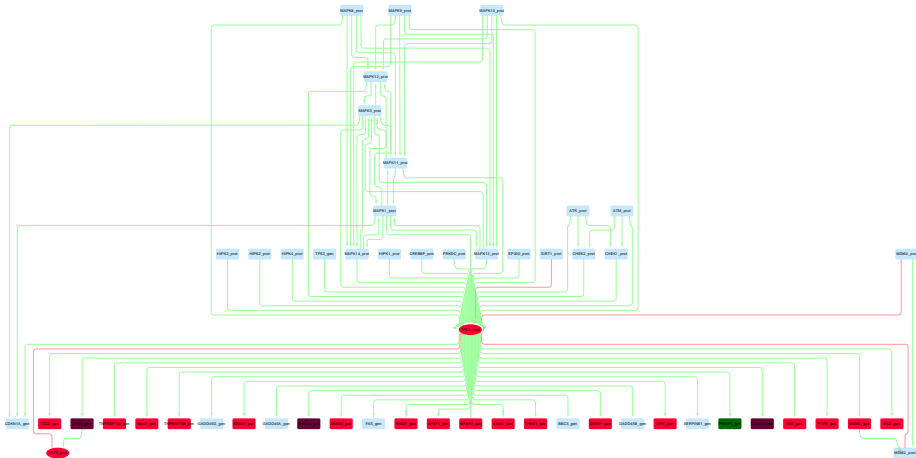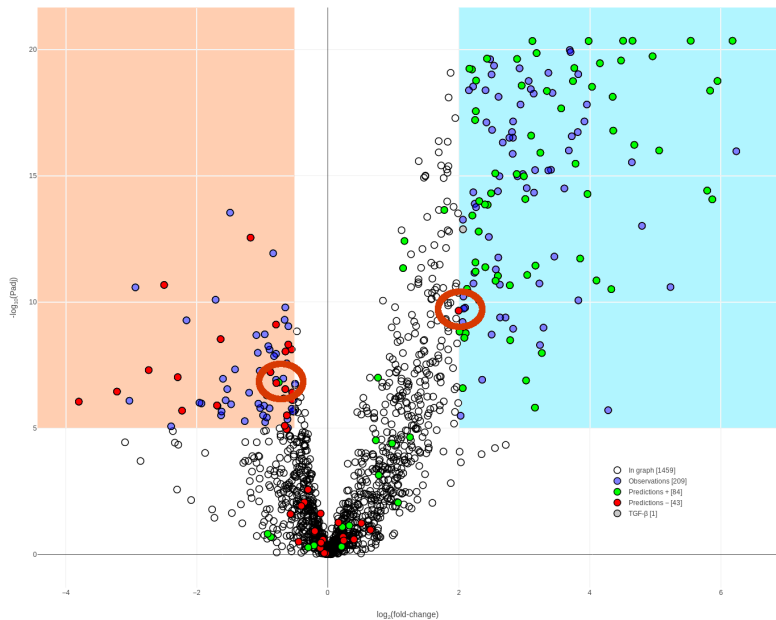# Hub example: TP53_prot



18 predictions directly depend of TP53_prot

Results of Iggy (predictions)

## Summary & Conclusion

### Summary

- Clustering + diff analysis: 2 lists of over- and under-expressed genes
- Graph extracted from Kegg: regulation + signaling
- 146 computational predictions (57 non-trivial)
- Predictions seem robust

### Objectives (to do)

- Explore survival curves compared to most robust genes
- Explore the literature regarding predicted complexes
  $\Rightarrow$ New proliferation signature?
- Try the same workflow on a different type of cancer (breast?)
- PUBLISH

# Bibliography I

📄 Cerami, E. G., Gross, B. E., Demir, E., Rodchenkov, I., Babur, Ö., Anwar, N., Schultz, N., Bader, G. D., and Sander, C. (2010).
Pathway Commons, a web resource for biological pathway data.
Nucleic acids research, 39.
http://www.pathwaycommons.org/.

📄 Hudson, T. J. and The International Cancer Genome Consortium (2010).
International network of cancer genome projects.
Nature, 464.
http://icgc.org/.

📄 Kanehisa, M., Furumichi, M., Tanabe, M., Sato, Y., and Morishima, K. (2017).
KEGG: new perspectives on genomes, pathways, diseases and drugs.
Nucleic Acids Research, 45(D1):D353–D361.
http://https://www.kegg.jp/.

📄 Lefebvre, M., Bourdon, J., Guziolowski, C., and Gaignard, A. (2017).
Regulatory and signaling network assembly through linked open data.
In Journées Ouvertes en Biologie, Informatique et Mathématiques.
Demo paper. https://github.com/symetric-group/bionets-demo.

Bibliography II

📄 Subramanian, A., Tamayo, P., Mootha, V. K., Mukherjee, S., Ebert, B. L., Gillette, M. A., Paulovich, A., Pomeroy, S. L., Golub, T. R., Lander, E. S., and Mesirov, J. P. (2005).
Gene Set Enrichment Analysis: A knowledge-based approach for interpreting genome-wide expression profiles.
Proc. of the Nat. Ac. of Sci., 102(43).
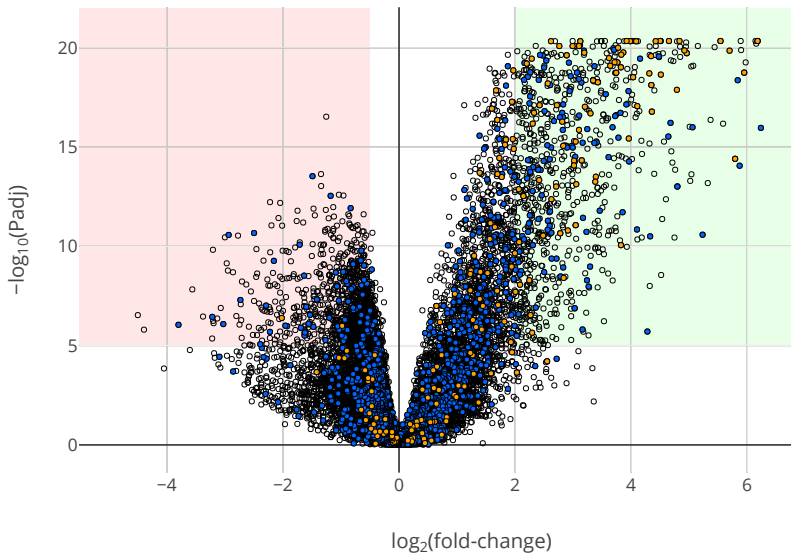http://software.broadinstitute.org/gsea/.

📄 Thiele, S., Cerone, L., Saez-Rodriguez, J., Siegel, A., Guziołowski, C., and Klamt, S. (2015).
Extended notions of sign consistency to relate experimental data to signaling and regulatory network topologies.
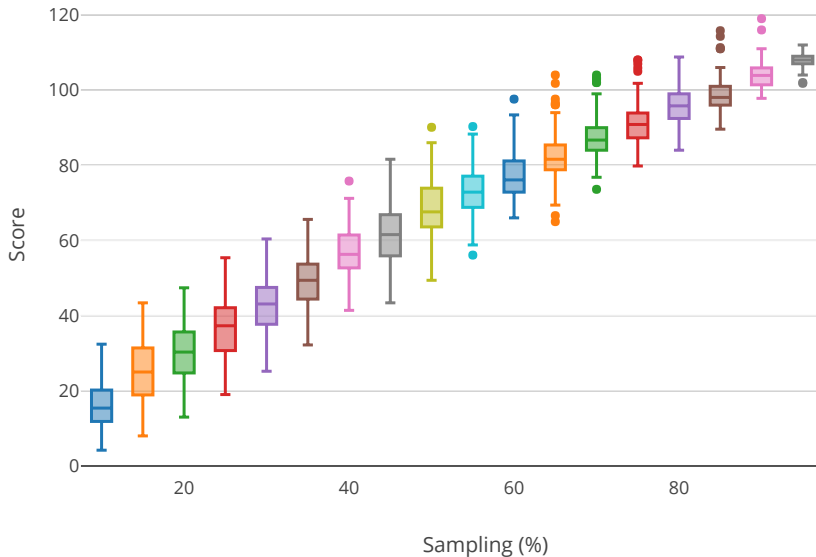BMC Bioinformatics, 16(1).
http://bioasp.github.io/iggy/.

# Initial ICGC data, EMT signature & genes found in Kegg

# Boxplot of the scores for each sampling

Evolution of max, min, mean and median of
good, bad and missing predictions compared to 100% sampling