



## Deconvolution and machine learning reveal oncogenic processes and characterize cancer patients

Petr Nazarov petr.nazarov@lih.lu

**LIH PI meeting** 



LUXEMBOURG INSTITUTE OF **HEALTH** RESEARCH DEDICATED TO LIFE

### FNR CORE Project and Project Proposals:



### $\textbf{DEMICS} \rightarrow \textbf{DIOMEDES} + \textbf{MEDEA}$



#### Work started: **2015** (with Laurent) Financed: **2018-2019** (with Francisco)





submitted 2021



Anna – Pl



#### submitted 2021



Michel – Co-PI

### 2018-2019



#### **Colleagues / co-authors**











# DEMICS

#### **Interns / students**









#### **Colleagues / collaborators**















Supported by Luxembourg National Research Fund C17/BM/11664971/DEMICS

### **Methods**





Hanahan D, Weinberg RA. Cell 2011, 144, 646-74



- Technical heterogeneity
- Native heterogeneity of biological tissues
- Inter/intra tumor heterogeneity due to clonal evolution

### **Methods**





### Methods





- + deterministic & fast
- + any number of samples
- + unsupervised
- often biological factors are presented by a sum of several components
- positive and negative values



- + correlates with biology
- + unsupervised (agnostic)
- + quite stable
- stochastic
- needs a lot of samples
- positive and negative values



- + semi-unsupervised
- + easy to interpret
- stochastic
- unstable

Sompairac, Nazarov, el al, Int J Mol Sci, 2019 (link) Cantini el al, Bioinformatics, 2019 (link)

### Consensus Independent Component Analysis (consICA)



#### 2021-05-07

LUXEMBOURG

INSTITUTE OF **HEALTH** 

### Melanoma







Collaboration with Dr S. Kreis (UL) Nazarov et al, BMC Medical Genomics, 2019 (link) In addition to diagnostics and prognostics, ICA allowed ranking patients based on activity of biological processes: cell cycle, signals of leukocytes, etc.

### **Melanoma**



#### **Deciphering biological processes and cell types**

Cluster	Compo- nent	Risk (p-value)	Meaning		P4PM	P6PM	P4NS	NHEM
Immune	RIC2	decreased (1.8e-4)	B cells		0.07	0.02	0.19	0.0
	RIC25	decreased (2.8e-7)	T cells		0.06	0.24	0.18	0.0
	RIC27	no effect	B cells	0.80	0.37	0.31	0.80	0.0
	RIC28	no effect	response to wounding	0.34	0.57	0.78	0.43	0.8
	RIC37	no effect	IFN signalling pathway	0.97	0.66	0.99	0.90	1.0
	RIC57	no effect	monocytes	0.00	0.25	0.24	0.02	0.0
	MIC20	decreased (1.2e-4)	T cells, chr1q32.2	0.14	0.08	0.37	0.02	0.1
Stromal and angiogenic	RIC13	no effect	cells of stroma	0.81	0.40	0.50	0.86	0.0
	RIC49	no effect	endothelial cells	0.73	0.12	0.29	0.84	0.0
	MIC22	no effect	miR-379/miR-410 cluster, chr14q32.2,14q32.31	0.29	0.20	0.27	0.38	0.1
	MIC25	no effect	stromal cells; clusters: chr1q24.3, 5q32, 17p13.1, 21q21.1	0.97	0.85	0.76	0.80	0.2
Skin-related	RIC5	increased (5.8e-3)	epidermis development and keratinisation	0.92	0.93	0.96	0.92	0.8
	RIC7	increased (8.9e-6)	epidermis development and keratinisation	0.94	0.93	0.93	0.95	0.5
	RIC19	increased (4.0e-2)	epidermis development and keratinisation	1.00	0.62	0.22	1.00	0.9
	RIC31	increased (2.2e-2)	epidermis development and keratinisation	0.98	0.85	0.89	0.99	0.2
	MIC9	increased (2.9e-2)	skin-specific miRNAs	0.95	0.88	0.87	0.91	0.8
Melanocytes	RIC4	increased (5.4e-3)	melanin biosynthesis	0.62	0.77	1.00	0.21	0.9
	RIC16	decreased (5.1e-4)	melanosomes (negative gene list)	0.68	0.77	0.54	0.75	0.3
	MIC11	no effect	potential regulators of malignant cells, chrXq27.3	0.21	0.96	0.62	0.13	0.4
	MIC14	decreased (1.5e-2)	potential regulators of melanocytes, chrXq26.3	0.01	0.29	0.67	0.29	0.3
Other	RIC55	increased (3.0e-2)	cell cycle	0.48	0.46	0.88	0.00	0.5
	RIC6	decreased (5.5e-3)	potentially linked to neuron differentiation	0.43	0.73	0.59	0.46	0.0
	MIC1	increased	regulators of EMT	0.11	0.07	0.02	0.19	0.0

#### ESTIMATE nature \_\_\_\_\_

#### Article | OPEN | Published: 11 October 2013

 $r^2 = 0.916$ 

score

mmune

Inferring tumour purity and stromal and immune cell admixture from expression

#### data

V

Kosuke Yoshihara, Maria Shahmoradgoli, Emmanuel Martínez, Rahulsimham Vegesna, Hoon Kim, Wandaliz Torres-Garcia, Victor Treviño, Hui Shen, Peter W. Laird, Douglas A. Levine, Scott L. Carter, Gad Getz, Katherine Stemke-Hale, Gordon B. Mills & Roel G.W. Verhaak



Nature Communications 4, Article number: 2612 (2013) Download Citation





Weights of RIC25

Weights of RIC13

#### Data integration: mRNA + miRNA + ...



 $\leftarrow$  New samples are mapped to the space defined by reference data.



#### Pancreatic cancers: ICA results of mRNA expression data from DKFZ cohort





N - healthy pancreas (41 samples)

**N.P** - histologically normal pancreas from patients with pancreatitis (15)

N.PDAC, N.TC, N.TE, N.TO - tumor-adjacent tissues (30+22+2+11)

P – pancreatitis (59)

- **PDAC** pancreas ductal adenocarcinomas (195)
- TC cystic tumors (24)
- **TE** neuroendocrine tumors (18)

**TO** – other tumors (31)

Components identified by ICA were annotated by biological functions (GO) and linked to survival using Cox regression.







Acc: 0.83	Ν	N.PDAC	Р	PDAC
pred.N	32	2.6	1.8	2
pred.N.PDAC	0.6	1.7	2	1.6
pred.P	4.7	17.3	51.8	5.4
pred.PDAC	3.7	8.4	3.4	186



≻

Aliaksandra Kakoichankava (MD student)

low cc

GERMAN CANCER RESEARCH CENTER IN THE HELMHOLTZ ASSOCIATION



PD Dr. med. Nathalia Giese

**Prof. Jörg Hoheisel** 

Dr. Andrea Bauer



### Gliomas → DIOMEDES + MEDEA

IDHwt

IDHmutnon-codel

IDHmut-

codel

⊙ n/a



#### **Gliomas in TCGA:**



#### (d) Histopathological differences







**IDHmut-codel** 

(c) deconvolved mRNA data

tSNE dimension 1

#### (e) ESTIMATE score and ICA (mRNA)



#### (f) Cell cycle component (mRNA)



#### **Cross-cohort prognosis**



50 8

### **GBM Cell Lines**





<u>.</u>

-0.08

-0.06

-0.04 -0.02

Gender

0.00

0.02

0.04

- We were able to map in-house cell line data onto TCGA dataset (GBM)
- Some components captured *technical factors* → (and thus clean other components from them)
- Other relevant biological information: cell cycle, cell migration, presence of stromal and immune cells. We were able to predict phenotype of cell lines using their transcriptomes.

### **GBM Cell Lines**





Patient-derived organoids and orthotopic xenografts of primary and recurrent gliomas represent relevant patient avatars for precision oncology

Anna Golebiewska<sup>1</sup> · Ann-Christin Hau<sup>1</sup> · Anaïs Oudin<sup>1</sup> · Daniel Stieber<sup>1,2</sup> · Yahaya A. Yabo<sup>1,3</sup> . Virginie Baus<sup>1</sup> · Vanessa Barthelemy<sup>1</sup> · Eliane Klein<sup>1</sup> · Sébastien Bougnaud<sup>1</sup> · Olivier Keunen<sup>1,4</sup> · May Wantz<sup>1</sup> · Alessandro Michelucci<sup>1,5,6</sup> · Virginie Neirinckx<sup>1</sup> · Arnaud Muller<sup>4</sup> · Tony Kaoma<sup>4</sup> · Petr V. Nazarov<sup>4</sup> · Francisco Azuaje<sup>4</sup> · Alfonso De Falco<sup>3,3,7</sup> · Ben Flies<sup>4</sup> · Lorraine Richart<sup>37,8,9</sup> · Suresh Poovathingal<sup>6</sup> · Thais Arns<sup>6</sup> · Kamil Grzyb<sup>6</sup> · Andreas Mock<sup>10,11,12,13</sup> · Christel Herold-Mende<sup>10</sup> · Anne Steino<sup>14,15</sup> · Dennis Brown<sup>14,15</sup> · Patrick May<sup>6</sup> · Hrvoje Miletic<sup>16,17</sup> · Tathiane M. Malta<sup>18</sup> · Houtan Noushmehr<sup>18</sup> · Yong-Jun Kwon<sup>7</sup> · Winnie Jahn<sup>19,20</sup> · Barbara Klink<sup>2,6,19,20,21</sup> · Georgette Tanner<sup>22</sup> · Lucy F. Stead<sup>22</sup> · Michel Mittelbronn<sup>6,7,8,9</sup> · Alexander Skupin<sup>6</sup> · Frank Hertel<sup>42,3</sup> · Rolf Bjerkviq<sup>1,16</sup> · Simone P. Niclou<sup>11,16</sup>



 ICA deconvolution is reasonable and predicts phenotypic behavior of cell lines

 Tumor cells show higher mobility in xenografts

#### **ESTIMATE** was confused



#### etr V. Nazarov<sup>4</sup>. Poornis Brown<sup>14,15</sup>. Dennis Brown<sup>14,15</sup>. Phenotype of cell lines were predicted using unsupervised deconvolution of their transcriptomes!

### **GBM Bulk-Sample Datasets** → **DIOMEDES**



#### IC39 – oxygen transport



⇒ consICA identified relevant & prognostic independent signals in a "homogenous" IDHwt datasets,

#### **IVY GAP**

**GLIOTRAIN** 

262 samples, 37 patients (~7 regions)

⇒ consICA characterize signals overrepresented in different tumor **niches** 



### **Pan-cancer & Multi-omics**





(MSc)

TCGA

The Cancer Genome Atlas

#### >11k patients, 33 types of tumors

- clinical data (age, gender, survival...)
- mRNA (10k samples, 20k features)
- miRNA (> 9k samples, ~1k features)
- methylation (>9k samples, 450k features)





Here we used consICA with 100 components & 40 runs

another example of ICA for methylation data:

Scherer M, Nazarov P, et al. Nature Protocols, 2020 (link)

#### protocols

#### PROTOCOI https://doi.org/10.1038/s41596-020-0369-

Check for updal

Reference-free deconvolution, visualization and interpretation of complex DNA methylation data using DecompPipeline, MeDeCom and FactorViz

Michael Scherer ©<sup>1,2</sup>, Petr V. Nazarov<sup>⊙3</sup>, Reka Toth<sup>©4,5</sup>, Shashwat Sahay<sup>1,7</sup>, Tony Kaoma<sup>3</sup>, Valentin Maurer<sup>4</sup>, Nikita Vedeneev<sup>6</sup>, Christoph Plass<sup>⊙4</sup>, Thomas Lengauer<sup>2</sup>, Jörn Walter<sup>⊙1</sup> and Pavlo Lutsik<sup>⊙4</sup>

### **Pan-cancer: ICA Components**



### **ICA Results: Cell Cycle**

#### **RIC27: Mitotic Cell Cycle**



Code 💌	Study Name	Ŧ		
ACC	Adrenocortical carcinoma			
BLCA	Bladder urothelial carcinoma			
BRCA	Breast invasive carcinoma			
CESC	Cervical sq. cell carcinoma and endocervical adenocarcinc			
CHOL	Cholangiocarcinoma			
COAD	Colon adenocarcinoma			
DLBC	Lymphoid neoplasm diffuse large b-cell lymphoma			
ESCA	Esophageal carcinoma			
GBM	Glioblastoma multiforme			
HNSC	Head and neck squamous cell carcinoma			
КІСН	Kidney chromophobe			
KIRC	Kidney renal clear cell carcinoma			
KIRP	Kidney renal papillary cell carcinoma			
LAML	Acute myeloid leukemia			
LCML	Chronic myelogenous leukemia			
LGG	Brain lower grade glioma			
LIHC	Liver hepatocellular carcinoma			
LUAD	Lung adenocarcinoma			
LUSC	Lung squamous cell carcinoma			
MESO	Mesothelioma			
ov	Ovarian serous cystadenocarcinoma			
PAAD	Pancreatic adenocarcinoma			
PCPG	Pheochromocytoma and paraganglioma			
PRAD	Prostate adenocarcinoma			
READ	Rectum adenocarcinoma			
SARC	Sarcoma			
SKCM	Skin cutaneous melanoma			
STAD	Stomach adenocarcinoma			
TGCT	Testicular germ cell tumors			
THCA	Thyroid carcinoma			
тнүм	Thymoma			
UCEC	Uterine corpus endometrial carcinoma			
UCS	Uterine carcinosarcoma			
UVM	Uveal melanoma			

### **Pan-cancer: ICA Components**



**RIC17: Signal of Mast Cells\*** 



Cox regression: logtest pv=1.4e-87 LHR=2.85 (CI = 2.57, 3.12)



Tumor-associated mast cells (TAMCs) ?

This can we wrong – we need a review from a biologist! ☺

(\*) assigned based on LM22 signature (CIBERSORT)

**RIC16: Signal of T-Cells\*** 



**RIC57: Angiogenesis** 



### **Pan-cancer: ICA-based Data Integration**





### **Pan-cancer: Classification**





**Pan-cancer: Prognosis** 







- Our modification of ICA deconvolution (*consICA*) :
  - Corrects technical biases
  - Extracts "cleaned" biological signals from bulk-sample data
  - > Maps new samples into the space of biologically meaningful components
  - Extracts prognostic features and features with classification power
  - Can be used to integrate multi-omics data
  - Diagnostic & prognostic properties could be expected for many cancers
- Was validated:
  - Using acceptable computational methods (cross-validation)
  - On cell lines
  - Independent cohorts of patients







### \* Single-cell part





Maryna Chepeleva (PhD student) (Bioinf.) **Tony** KAOMA (Bioinf.)



Anna GOLEBIEWSKA (PI)



Simone NICLOU

### ICA for Single Cell RNA-seq Data



### **Correction of technical effects**



$$E_{nm} \Rightarrow S_{nk} \times M_{km}$$
$$M'_{7,m} \leftarrow 0$$
$$E'_{nm} \leftarrow S_{nk} \times M'_{km}$$

t-SNE representations of original data (A) and ICA-recovered data, after excluding batch effect (B) or several (C) components linked to technical factors ("cell size").



#### ARTICLE

https://doi.org/10.1038/s41467-019-09853-z

Stem cell-associated heterogeneity in Glioblastoma results from intrinsic tumor plasticity shaped by the microenvironment

Anne Dirkse<sup>1,2,12</sup>, Anna Golebiewska<sup>®</sup><sup>1,12</sup>, Thomas Buder<sup>3,4</sup>, Petr V. Nazarov<sup>®</sup><sup>5</sup>, Arnaud Muller<sup>5</sup>, Suresh Poovathingal<sup>6</sup>, Nicolaas H.C. Brons<sup>7</sup>, Sonia Leite<sup>8</sup>, Nicolas Sauvageot<sup>8</sup>, Dzjemma Sarkisjan<sup>1</sup>, Mathieu Seyfrid<sup>1</sup>, Sabrina Fritah<sup>1</sup>, Daniel Stieber<sup>1</sup>, Alessandro Michelucci<sup>®</sup><sup>1,6</sup>, Frank Hertel<sup>9</sup>, Christel Herold-Mende<sup>10</sup>, Francisco Azuaje<sup>®</sup>, Alexander Skupin<sup>6</sup>, Rolf Bjerkvig<sup>1,1</sup>, Andreas Deutsch<sup>3</sup>, Anja Voss-Böhme<sup>3,4</sup> & Simone P. Niclou<sup>®</sup><sup>1</sup>

### ICA for Single Cell RNA-seq Data

LUXEMBOURG INSTITUTE OF **HEALTH** RESEARCH DEDICATED TO LIFE

### **Cell Cycle in Single Cells**







Dominiguez (2016) Cell Research



Dirkse, Golebiewska et al. Nature Communications, 2019 (link) Sompairac, Nazarov el al. Int J Mol Sci, 2019 (link)

### **Combining Bulk and Single Cell**

LUXEMBOURG INSTITUTE OF **HEALTH** RESEARCH DEDICATED TO LIFE



### **GBM: Single Cell Data**





### **GBM: Single Cell Data Deconvolution**





*consICA* can assess cellular subpopulations and phenotypic states associated with specific biological processes

### Longitudinal PDOX: TMZ Resistance





Examples of independent components that represent tumor variability:

- between patients (left)
- within each patient regardless treatment (middle)
- within each patient before and after treatment (right)

Deconvolution captured signal related to treatment (no genes using standard DEA approach) GO:BP

translation initiation mRNA catabolic process cytoplasmic translation **GO:CC** cytosolic ribosome **GO:MF** structural constituent of ribosome RNA binding



- Deconvolution with *consICA* on single cell data:
  - Corrects technical biases
  - Extracts signals of biological processes (or cell types) from single-cell data
  - Could detect weak signals that are masked by other processes (e.g. TMZ resistance masked by cell cycle and inter-tumor variability)
  - Can be used to interpret results of bulk-sample data deconvolution





# **MEDEA**





Olivier

Sang Yoon KEUNEN



Michel MITTELBRONN (Co-PI)

2021-05-07

ΚΙΜ

### Background







Native heterogeneity of tissues

Inter/intra tumor heterogeneity

Issues in histopathological image analysis:

- Tedious analysis
- In some cancers (e.g. prostate) < 1% of the image is cancer-related
- Standard approaches require supervised
  "pixel-wise" labelling unrealistic



### **Observation 1**





International Journal of *Molecular Sciences* 



LUXEMBOURG

MDPI

#### Review Independent Component Analysis for Unraveling the Complexity of Cancer Omics Datasets

Nicolas Sompairac <sup>1,2,3,4</sup>, Petr V. Nazarov <sup>5</sup>, Urszula Czerwinska <sup>1,2,3</sup>, Laura Cantini <sup>6</sup>, Anne Biton <sup>7</sup>, Askhat Molkenov <sup>8</sup>, Zhaxybay Zhumadilov <sup>8,9</sup>, Emmanuel Barillot <sup>1,2,3</sup>, Francois Radvanyi <sup>1,10</sup>, Alexander Gorban <sup>11,12</sup>, Ulykbek Kairov <sup>8</sup>, and Andrei Zinovyev <sup>1,2,3,\*</sup>

### Observation

ICA results are confirmed by H&E histopathology: smooth muscles, fibroblasts, cell cycle were observed.

### **Observation 2**



#### ICA results of mRNA expression data from TCGA-PAAD cohort



### **Idea for Approach**



#### Patient classification using weakly supervised DLN

ARTICLES https://doi.org/10.1038/s41591-019-0508-1

### Clinical-grade computational pathology using weakly supervised deep learning on whole slide images

Gabriele Campanella<sup>1,2</sup>, Matthew G. Hanna<sup>1</sup>, Luke Geneslaw<sup>1</sup>, Allen Miraflor<sup>1</sup>, Vitor Werneck Krauss Silva<sup>1</sup>, Klaus J. Busam<sup>1</sup>, Edi Brogi<sup>1</sup>, Victor E. Reuter<sup>1</sup>, David S. Klimstra<sup>1</sup> and Thomas J. Fuchs<sup>1</sup>,<sup>2,\*</sup>

Deep Convolutional Network (DCN)



#### **RNA-seq prediction**



ARTICLE

#### https://doi.org/10.1038/s41467-020-17678-4 OPEN



## A deep learning model to predict RNA-Seq expression of tumours from whole slide images

Benoît Schmauch <sup>1</sup><sup>™</sup>, Alberto Romagnoni<sup>1,4</sup>, Elodie Pronier<sup>1,4</sup>, Charlie Saillard<sup>1</sup>, Pascale Maillé<sup>2,3</sup>, Julien Calderaro<sup>2,3</sup>, Aurélie Kamoun <sup>1</sup>, Meriem Sefta<sup>1</sup>, Sylvain Toldo<sup>1</sup>, Mikhail Zaslavskiy<sup>1</sup>, Thomas Clozel <sup>1</sup>, Matahi Moarii<sup>1</sup>, Pierre Courtiol<sup>1,5</sup> & Gilles Wainrib<sup>1,5™</sup>

### Idea

Instead of predicting mixed, bulk-sample mRNA or DNA-methylation signals, we will predict already deconvolved, clean signals, linked to biological processes & cell subpopulations.

medicine

### Artificial Neural Networks → Deep Learning Networks

LUXEMBOURG INSTITUTE OF **HEALTH** RESEARCH DEDICATED TO LIFE



Multilayer perceptron, a.k.a. (Deep) feed-forward network, back-propagation network fully-connected layers, etc...



My first "love"... 🙂

Nazarov et al (2004) J Chem Inf Comput Sci







MLP

### **MEDEA: Project Overview**





CAE: convolutional autoencoder; CNN: convolutional neural network; FC: fully-connected network or layer; ICA: independent component analysis; ML: machine learning; ROI: region of interest; WSI: whole slide image.

(a) Deconvolution of the omics data using developed tool *consICA*. This method was already developed and applied to entire GTEx (mRNA), TCGA (mRNA and meDNA), and DKFZ (mRNA) cohorts.

**(b) Image analysis and feature extraction** starts with a pretrained *Xception* model and uses weakly supervised training to fine-tune model's parameters. Two strategies will be compared in the project: strategy 1 is a semi-supervised one using CNNbased classifier and strategy 2 – completely unsupervised using CAE. *Xception* will be used as an initial estimation of the encoder's parameters.

(c) Integration of ICA-weights and image features will be done either by a classical ML-approach (linear regression or random forest regression) or by a FC neural network.

(d) A thorough validation of the results include (i) validation of an external pancreatic cancer cohort (DKFZ) and collection and (ii) in-depth analysis of in-house (LNS) samples of glioma patients. The expertise of the Co-PI (pathologist) will be used to validated predictions and the PI and his team will control that the WSI-features are sensible and not artefacts.

### **Preliminary Results**





LIH PI Meeting 37

### **Tile-level Feature Extraction**





These features were summarized to slide-level. Only 50% topcorrelated tiles were preserved (can be further improved later...)

### **Slide-level Analysis and ICA**

![](_page_38_Figure_1.jpeg)

![](_page_38_Figure_2.jpeg)

### **Predictions**

![](_page_39_Picture_1.jpeg)

FDR

8e-28

8e-28

4e-22

6e-22

7e-20

![](_page_39_Figure_2.jpeg)

WSI-features

### Predicting genes

- 0.4% of the genes showed  $R^2>0.9$
- $28\% R^2 > 0.5$

### **Predicting ICs or Genes?**

![](_page_40_Picture_1.jpeg)

#### **Two best-predicted coding genes**

![](_page_40_Figure_3.jpeg)

#### GO:BP pos : 32 terms(FDR<0.01)

Term	FDR
neutrophil degranulation	1.58e-27
defense response to other organism	4.36e-20
immune response	1.48e-15
neutrophil chemotaxis	8.71e-12

![](_page_40_Figure_6.jpeg)

# Majority of the predicted genes: tissue-specific non-coding

### No GO enriched!

LIH PI Meeting 41

![](_page_41_Picture_1.jpeg)

- > Deep Learning Networks could be used for feature extraction
- Image features could be used to predict deconvolved signals
- Deconvolved ("clean") signals are better predicted than genes (and related GO gene sets)
- > Approach was validated on a part of GTEx data