# Combination of Multi-modal Data for Improved Patient Characterization
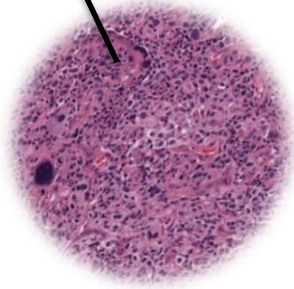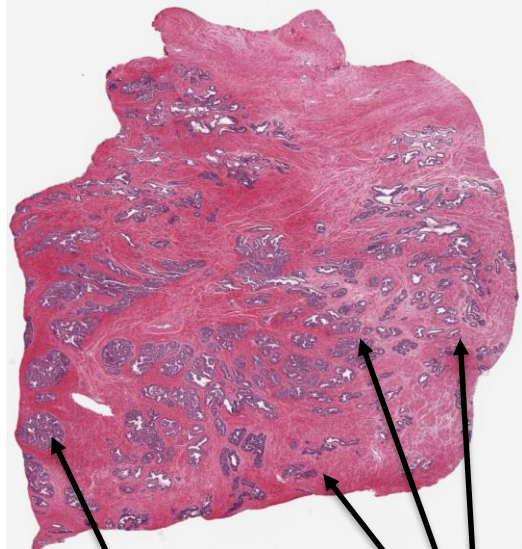
## Petr Nazarov

# Heterogeneity

## Levels of Heterogeneity in Samples of Cancer Patients



Native h. of biological tissues

Normal cells | Immune cells | Fibroblasts

Cancer cells

Invasive cancer cells

Hanahan, Weinberg. *Cell* **2011**, 144, 646-74

Technical heterogeneity

Cancer sample → Data

Inter- and intra-tumor heterogeneity

Clonal evolution

Cancer stem-cells

Cell plasticity

Neftel, et al. *Cell* **2019**, 178:835

Dirkse, et al. *Nat Commun* **2019**, 10:1787

Tirosh, et al. *Science* **2016**, 352(6282):189

De Sousa E Melo, et al. *EMBO* **2013**, 14(8):686

# Invasive Approach 1: Histopathology

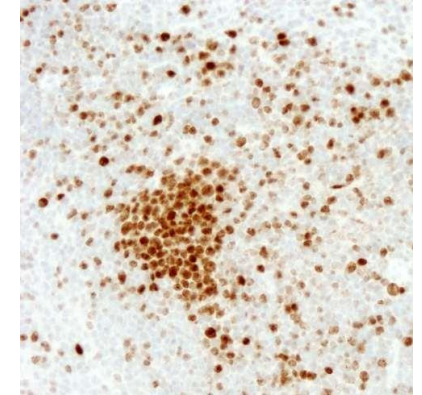**Hematoxylin and Eosin (H&E) stain**



Tumor: 1%  Normal: 99%

**Features of histopathology**

➢ Gold standard!

➢ Cheap (H&E or 2-3 antibodies in IHC)

➢ Captures native heterogeneity of tissues

➢ Shows inter/intra tumor heterogeneity

➢ Often allows precise diagnostics
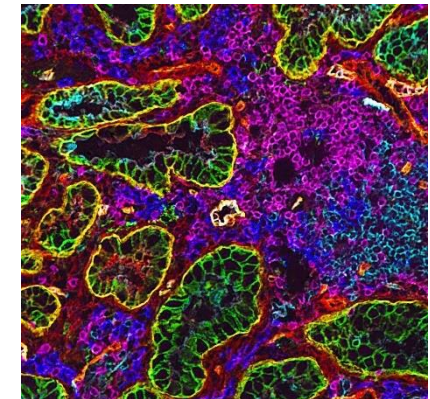
**Issues in histopathological image analysis:**

➢ Tedious analysis

➢ In some cancers (e.g. prostate) < 1% of the image is cancer-related

➢ For some cancers, it does not allow precise diagnostics (e.g. some astrocytomas vs oligodendrogliomas)

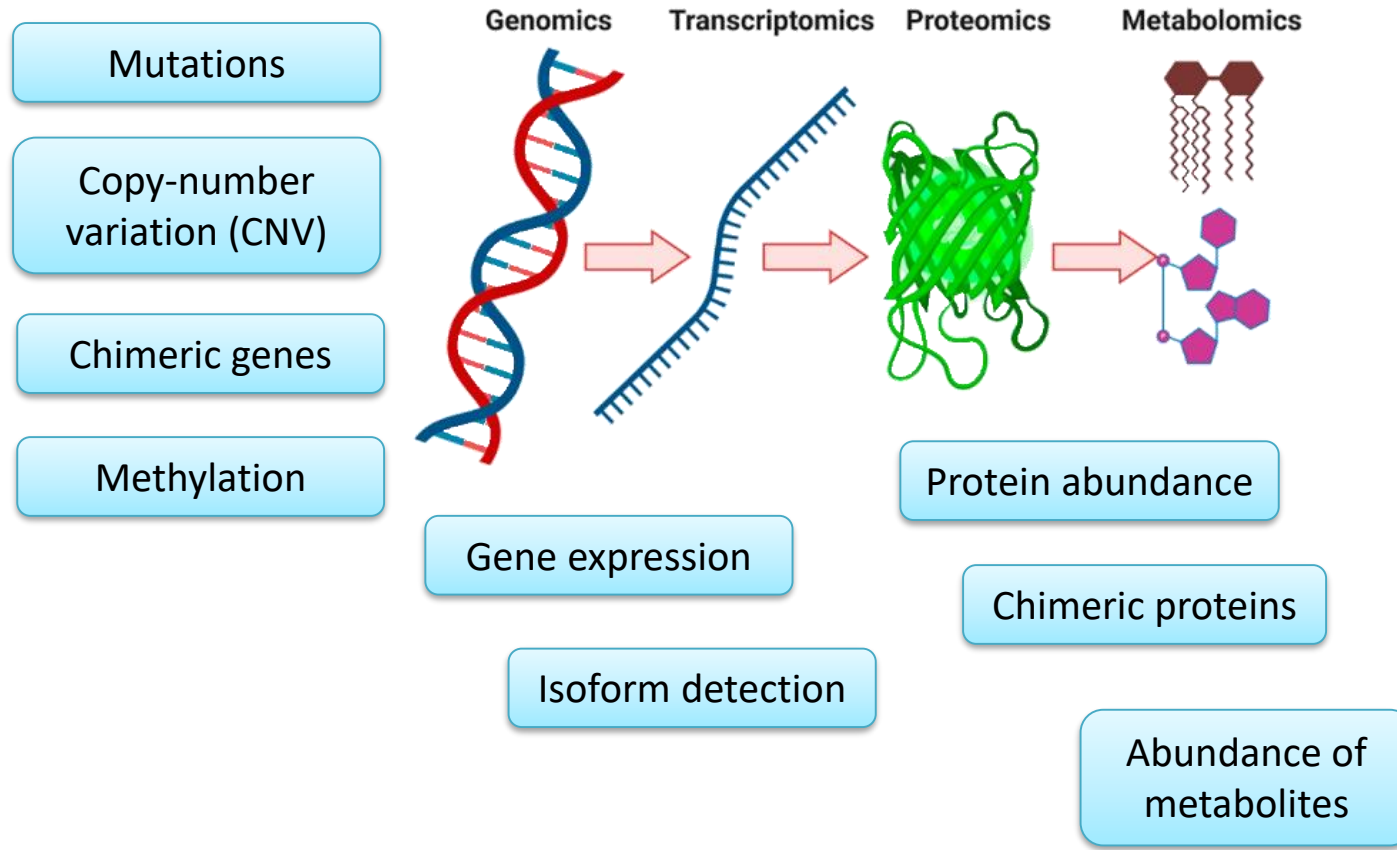➢ Gives non-structured data

**Immunohistochemistry (IHC)**



*Ki-67 - proliferation marker*

**Multicolor IHC**

# Invasive Approach 2: Molecular Profiling



**Features of molecular approach**

➢ Very specific

➢ Generate a lot of data

➢ Generate structured data

**Issues of molecular approach**

➢ Quite expensive

➢ Is sensitive to heterogeneity of samples

➢ Is sensitive to a technique

Mutations

Copy-number variation (CNV)

Chimeric genes

Methylation

Gene expression

Isoform detection

Protein abundance

Chimeric proteins

Abundance of metabolites

## DNA methylation–based classification of central nervous system tumours

Capper et al. *Nature* **2018**, 555(7697):469
Capper et al. *Acta Neuropathologica* **2018**, 136:181

➤ Methylation showed more specificity than histopathology identifying types of brain tumors

➤ Highly standardized pipeline allowed analysis across many cohorts

➤ **Result**: "Heidelberg classifier" is used by pathologists ☺

# Improvements

**1. Histopathology**



**2. Molecular methods**



➢ Automate analysis
➢ Transform unstructured data (images) to structured (features)

➢ Deconvolute mixed signals
➢ Integrate various molecular data

Integrate both approaches for better patient diagnostics and studying molecular processes

# 1. Digital Histopathology and Feature Extraction



**The Task**

N slides:
$10^4 \times 10^5$ pixels

N x M tiles / patches:
256 x 256 pixels

**1 patient**

AI → *Type*

**1 label**

~ $10^3$-$10^4$ tiles

**Classical image analysis approaches**

Counting nuclei
Edge selection
Cell shape
Cell graph
....

→ features

**Deep Artificial Neural Networks**

Deep convolutional neural network (CNN)

class

features

Convolutional Autoencoder (CAE)

profile $S_1$

$X_1 = S_1 + S_2$

$X_3 = S_1 + 3S_2$

Mixing

$X = S \times M$

$M =$

| 1 | 2 | 1 |
|---|---|---|
| 1 | 1 | 3 |

profile $S_2$

molecular features

De-mixing

$X_2 = 2S_1 + S_2$

Often called:
- **decomposition**
- **deconvolution**

# Deconvolution via Matrix Factorization

**Data**

genes

$E_{nm}$

samples

$\approx$

**Signals / Pattern**

genes

$S_{nk}$

$\times$

components

**Weights / Amplitude**

$M_{km}$

samples

components

**Matrix tri-factorization**

$$A_1, A_2, A_3 \approx G * S_1, S_2, S_3 * G^T$$

$$\min_{G \geq 0, S} \sum_{i=1}^{3} \| A_i - G \cdot S_i \cdot G^T \|_F^2$$

Malod-Dognin et al. *Nat Commun* **2019**, 10:805

**Multi-omics Factor Analysis**

Argelaguet et al. *Mol Syst Biol* **2018**, 14:e8124

**PCA**: principal component analysis
**NMF**: non-negative matrix factorization
**ICA**: independent component analysis
*etc.*

# Deconvolution Methods

## PCA



deterministic

variability

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

## ICA



stochastic

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

## NMF



stochastic

cells B

cells A

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

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

**Joined Expression Data**

**Independent Signals**

genes / samples

$E_{nm}$

**ICA**

multiple runs

$S_{nk}$ × $M_{km}$

**Weights**

samples / components

**Diagnostics:** using machine learning tools to predict classes of the samples

Weights *M* in patient groups

weights / patient groups

*M*

**Prognostics:** using Cox regression & combine weights into a risk score $RS_j$ to patient survival

$$RS_j = \sum_{i=1}^{i=k} R_i^2 H_i M_{i,j}^*$$

**Functional annotation:** linking components to biological processes and cell types

Genes, contributing to one component

contribution / genes (ordered)

*S*

**Data driven!**

**Discovery dataset (TCGA)**

**Investigation dataset (new patients)**

Nazarov et al. BMC Medical Genomics (2019) 12:132
https://doi.org/10.1186/s12920-019-0578-4

BMC Medical Genomics

TECHNICAL ADVANCE — Open Access

Deconvolution of transcriptomes and miRNomes by independent component analysis provides insights into biological processes and clinical outcomes of melanoma patients

Petr V. Nazarov, Anke K. Wienecke-Baldacchino, Andrei Zinovyev, Urszula Czerwińska, Arnaud Muller, Dorothée Nashan, Gunnar Dittmar, Francisco Azuaje and Stephanie Kreis

**consICA:** Nazarov et al **BMC Medical Genomics**, 2019 (link)
ICA review: Sompairac, et al **Int J Mol Sci**, 2019 (link)
Application: Golebiewska et al, **Acta Neuropathol**, 2020
Scherer, Nazarov et al, **Nat Protoc**, 2020

**Reference dataset**
530 GBM patients
(TCGA)

**Investigated dataset**
58 samples:
cell lines, xenografts &
patient tissues

**ICA**

**Biological knowledge:**
bio-processes
and sample
composition

**Technical/trivial components: gender and platforms**



➤ 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

**ICA correctly predicts sample composition & phenotype**



Acta Neuropathologica (2020) 140:919–949
https://doi.org/10.1007/s00401-020-02226-7

ORIGINAL PAPER

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

Anna Golebiewska[1] · Ann-Christin Hau[1] · Anaïs Oudin[1] · Daniel Stieber[1,2] · Yahaya A. Yabo[1,3] · Virginie Baus[1] · Vanessa Barthelemy[1] · Eliane Klein[1] · Sébastien Bougnaud[1] 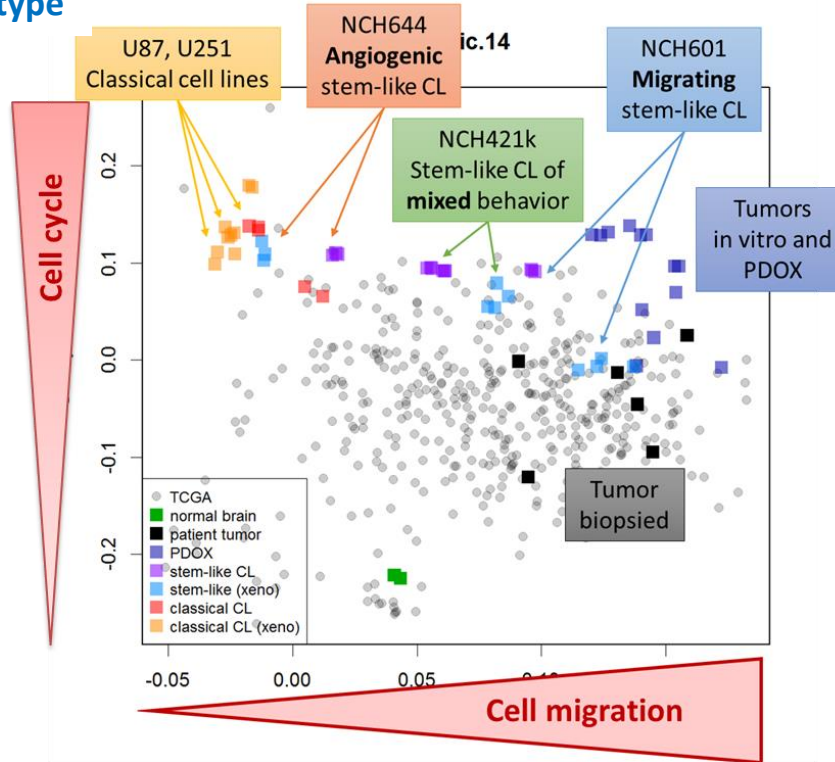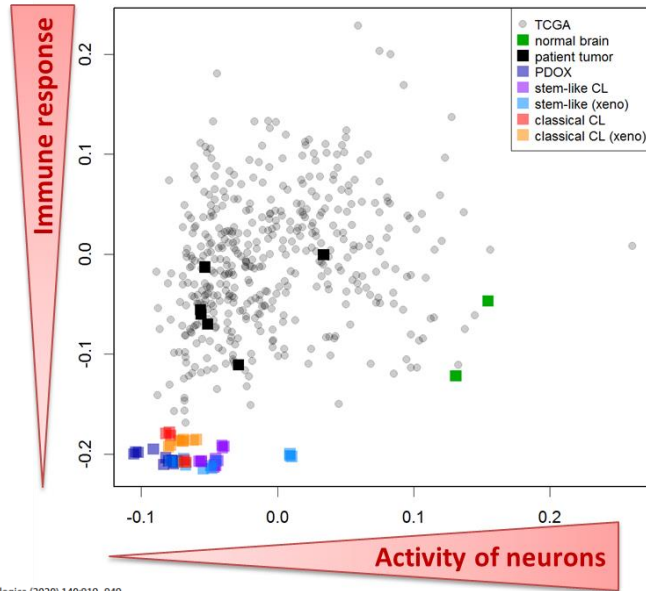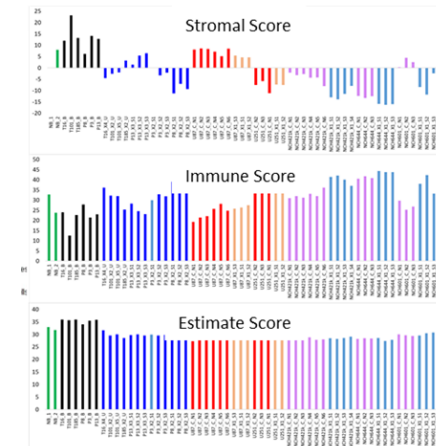· Olivier Keunen[1,4] · May Wantz[1] · Alessandro Michelucci[1,5,6] · Virginie Neirinckx[1] · Arnaud Muller[4] · Tony Kaoma[4] · Petr V. Nazarov[4] · Francisco Azuaje[4] · Alfonso De Falco[2,3,7] · Ben Flies[2] · Lorraine Richart[3,7,8,9] · Suresh Poovathingal[6] · Thais Arns[6] · Kamil Grzyb[6] · Andreas Mock[10,11,12,13] · Christel Herold-Mende[10] · Anne Steino[14,15] · Dennis Brown[14,15] · Patrick May[6] · Hrvoje Miletic[16,17] · Tathiane M. Malta[18] · Houtan Noushmehr[18] · Yong-Jun Kwon[9] · Winnie Jahn[19,20] · Barbara Klink[2,9,19,20,21] · Georgette Tanner[22] · Lucy F. Stead[22] · Michel Mittelbronn[6,7,8,9] · Alexander Skupin[6] · Frank Hertel[6,23] · Rolf Bjerkvig[1,16] · Simone P. Niclou[1,16]

- ➢ ICA deconvolution is reasonable and predicts phenotypic behavior of cell lines

- ➢ Tumor cells show higher mobility in xenografts

**ESTIMATE was confused**



Golebiewska A. et al, **Acta Neuropathologica, 2020** (link)

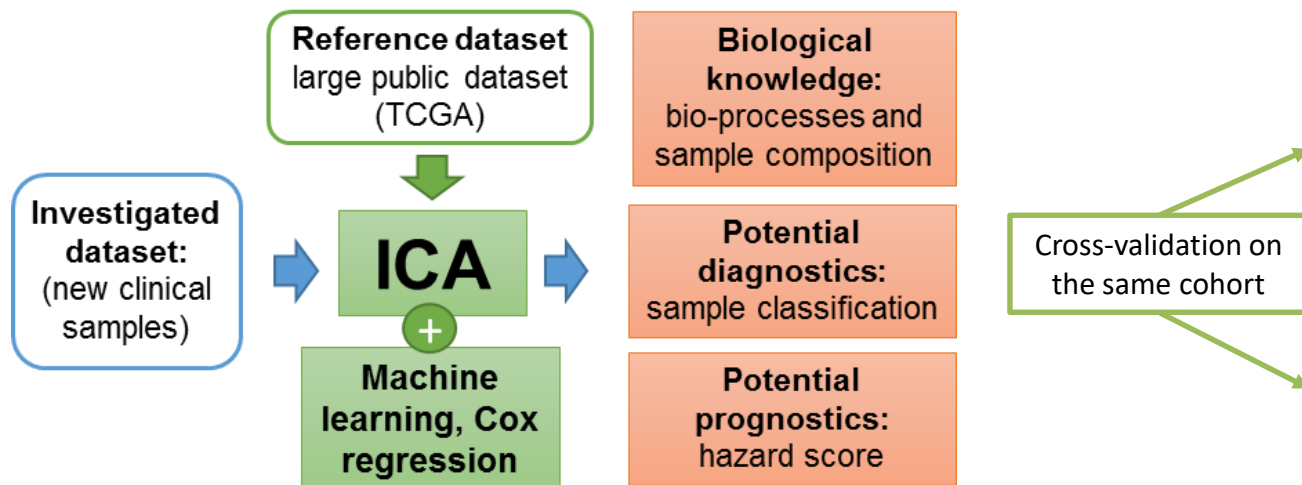Phenotype of cell lines were predicted using unsupervised deconvolution of their transcriptomes!

**Reference dataset**
large public dataset
(TCGA)

**Investigated dataset:**
(new clinical samples)

**ICA** +

**Machine learning, Cox regression**

**Biological knowledge:**
bio-processes and sample composition

**Potential diagnostics:**
sample classification

**Potential prognostics:**
hazard score

Cross-validation on the same cohort

| Cluster | | | |
|---|---|---|---|
| Accuracy | **Actual cluster** | | |
| 90.0% | **immune** | **keratine** | **MITF-low** |
| **immune** | 160 | 9 | 6 |
| **keratine** | 9 | 91 | 6 |
| **MITF-low** | 1 | 2 | 47 |

$$RS_j = \sum_{i=1}^{i=k} R_i^2 H_i M_{i,j}^*$$

$j$ – patient index
$i$ – component index
$R^2_i$ – stability of $i$-th component (from 0 to 1)
$H_i$ – Cox' log hazard ratio calculated on **training set**
$M^*_{i,j}$ – element of centered & scaled M-matrix



A  Discovery cohort
P-value (log-rank)= 5.6e-16
LHR= 0.49 (CI = 0.37, 0.61)

B  Validation cohort
P-value (log-rank) = 1.3e-03
LHR= 0.87 (CI = 0.28, 1.45)

Independent cohort, different platform

➢ In addition to diagnostics and prognostics, ICA allowed ranking patients based on the activity of biological processes: cell cycle, signals of leukocytes, etc.

Nazarov et al, BMC Medical Genomics, 2019

## Deciphering biological processes and cell types

| Cluster | Component | Risk (p-value) | Meaning | P2PM | P4PM | P6PM | P4NS | NHEM |
|---|---|---|---|---|---|---|---|---|
| Immune | RIC2 | decreased (1.8e-4) | B cells | 0.11 | 0.07 | 0.02 | 0.19 | 0.01 |
| | RIC25 | decreased (2.8e-7) | T cells | 0.26 | 0.06 | 0.24 | 0.18 | 0.00 |
| | RIC27 | no effect | B cells | 0.80 | 0.37 | 0.31 | 0.80 | 0.00 |
| | RIC28 | no effect | response to wounding | 0.34 | 0.57 | 0.78 | 0.43 | 0.84 |
| | RIC37 | no effect | IFN signalling pathway | 0.97 | 0.66 | 0.99 | 0.90 | 1.00 |
| | RIC57 | no effect | monocytes | 0.00 | 0.25 | 0.24 | 0.02 | 0.00 |
| | MIC20 | decreased (1.2e-4) | T cells, chr1q32.2 | 0.14 | 0.08 | 0.37 | 0.02 | 0.19 |
| Stromal and angiogenic | RIC13 | no effect | cells of stroma | 0.81 | 0.40 | 0.50 | 0.86 | 0.03 |
| | RIC49 | no effect | endothelial cells | 0.73 | 0.12 | 0.29 | 0.84 | 0.00 |
| | MIC22 | no effect | miR-379/miR-410 cluster, chr14q32.2,14q32.31 | 0.29 | 0.20 | 0.27 | 0.38 | 0.16 |
| | MIC25 | no effect | stromal cells; clusters: chr1q24.3, 5q32, 17p13.1, 21q21.1 | 0.97 | 0.85 | 0.76 | 0.80 | 0.26 |
| Skin-related | RIC5 | increased (5.8e-3) | epidermis development and keratinisation | 0.92 | 0.93 | 0.96 | 0.92 | 0.87 |
| | RIC7 | increased (8.9e-6) | epidermis development and keratinisation | 0.94 | 0.93 | 0.93 | 0.95 | 0.57 |
| | RIC19 | increased (4.0e-2) | epidermis development and keratinisation | 1.00 | 0.62 | 0.22 | 1.00 | 0.93 |
| | RIC31 | increased (2.2e-2) | epidermis development and keratinisation | 0.98 | 0.85 | 0.89 | 0.99 | 0.28 |
| | MIC9 | increased (2.9e-2) | skin-specific miRNAs | 0.95 | 0.88 | 0.87 | 0.91 | 0.83 |
| Melanocytes | RIC4 | increased (5.4e-3) | melanin biosynthesis | 0.62 | 0.77 | 1.00 | 0.21 | 0.96 |
| | RIC16 | decreased (5.1e-4) | melanosomes (negative gene list) | 0.68 | 0.77 | 0.54 | 0.75 | 0.39 |
| | MIC11 | no effect | potential regulators of malignant cells, chrXq27.3 | 0.21 | 0.96 | 0.62 | 0.13 | 0.48 |
| | MIC14 | decreased (1.5e-2) | potential regulators of melanocytes, chrXq26.3 | 0.01 | 0.29 | 0.67 | 0.29 | 0.38 |
| Other | RIC55 | increased (3.0e-2) | cell cycle | 0.48 | 0.46 | 0.88 | 0.00 | 0.53 |
| | RIC6 | decreased (5.5e-3) | potentially linked to neuron differentiation | 0.43 | 0.73 | 0.59 | 0.46 | 0.01 |
| | MIC1 | increased (9.4e-4) | regulators of EMT | 0.11 | 0.07 | 0.02 | 0.19 | 0.01 |

ESTIMATE



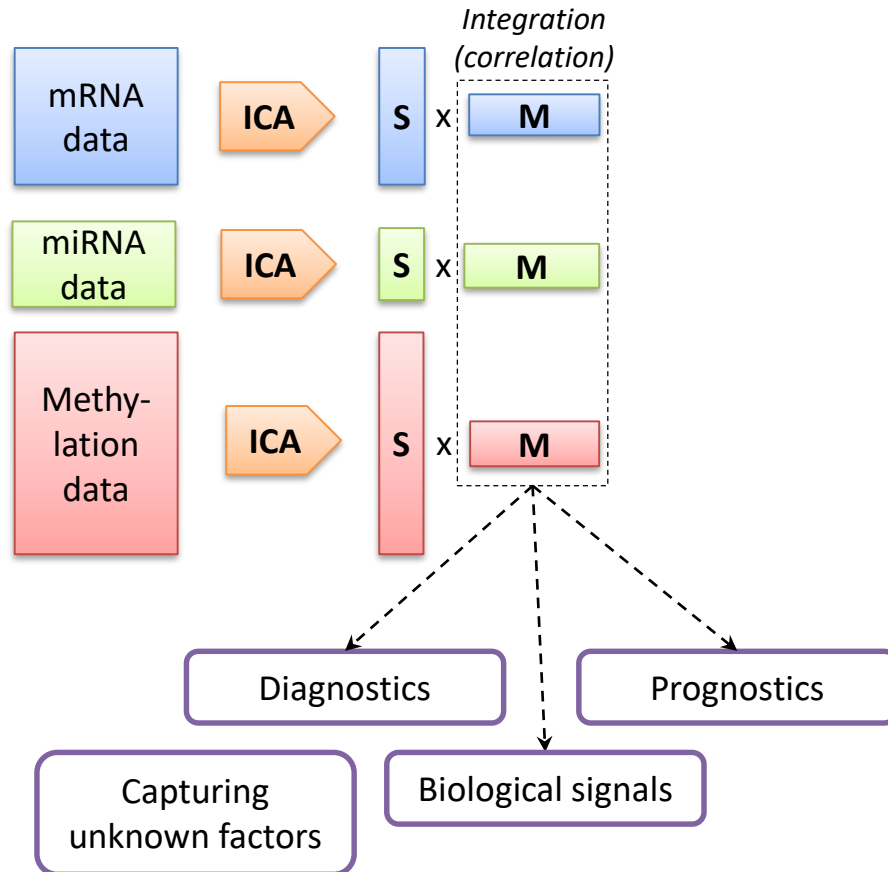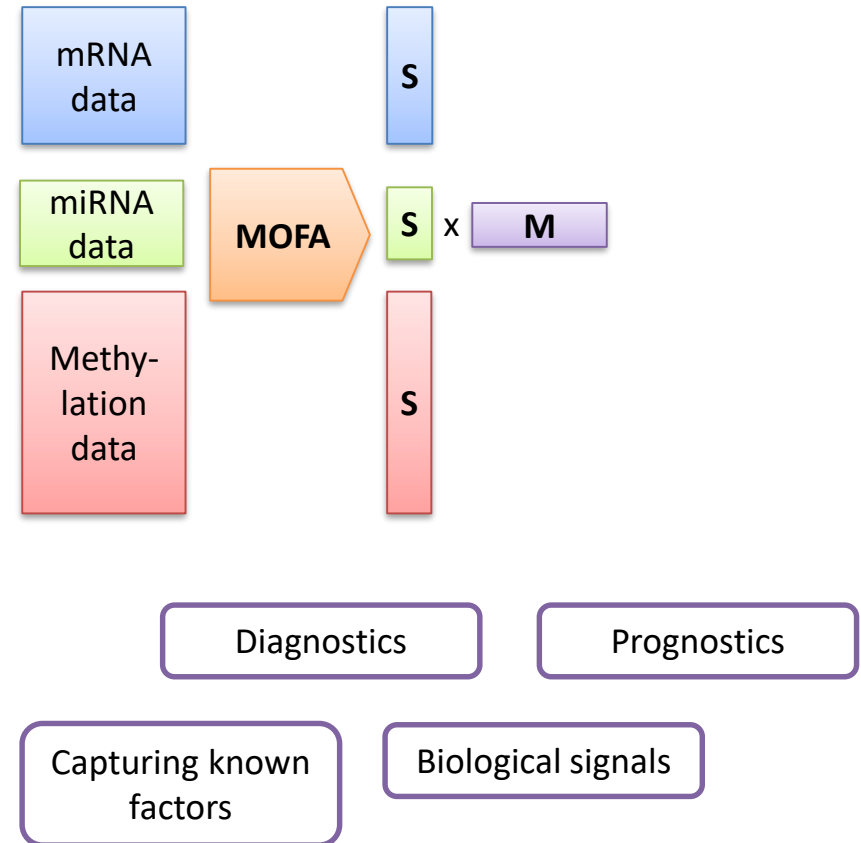## Data integration: mRNA + miRNA + …



← New samples are mapped to the space defined by reference data.

# Multi-omics Data Integration via Deconvolution
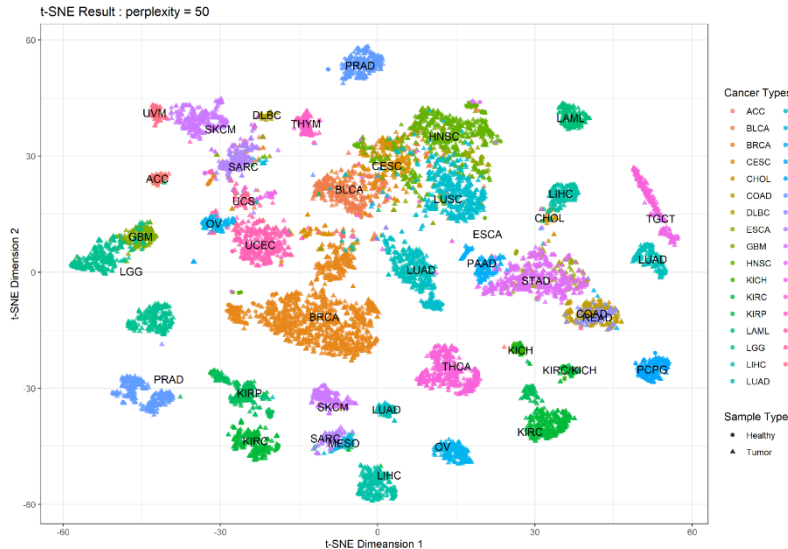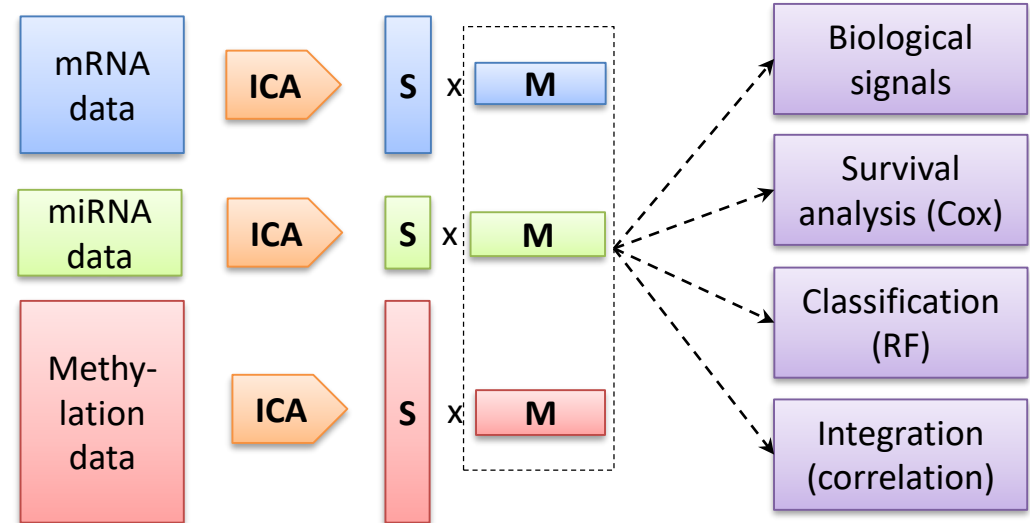
# Pan-Cancer Data Integration

## 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)
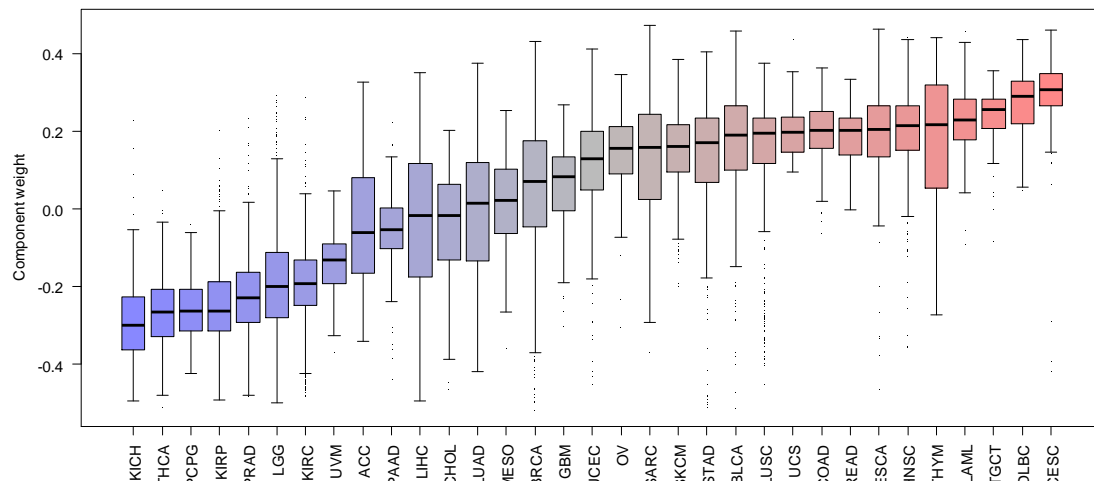- **methylation** (>9k samples, 450k features)

### Approach



Here we used *consICA* with 100 components & 40 runs

## ICA Results: Cell Cycle

### RIC27: Mitotic Cell Cycle

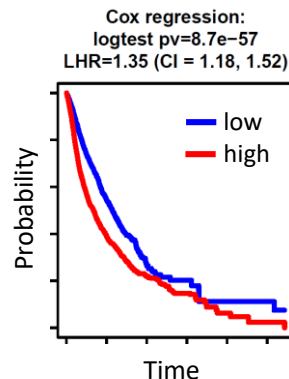

Cox regression:
logtest pv=8.7e−57
LHR=1.35 (CI = 1.18, 1.52)

prostate adenocarcinoma

low grade glioma

glioblastoma (IDHwt)

diffuse lymphoma

cervical s.c.c & endoservical a.c.
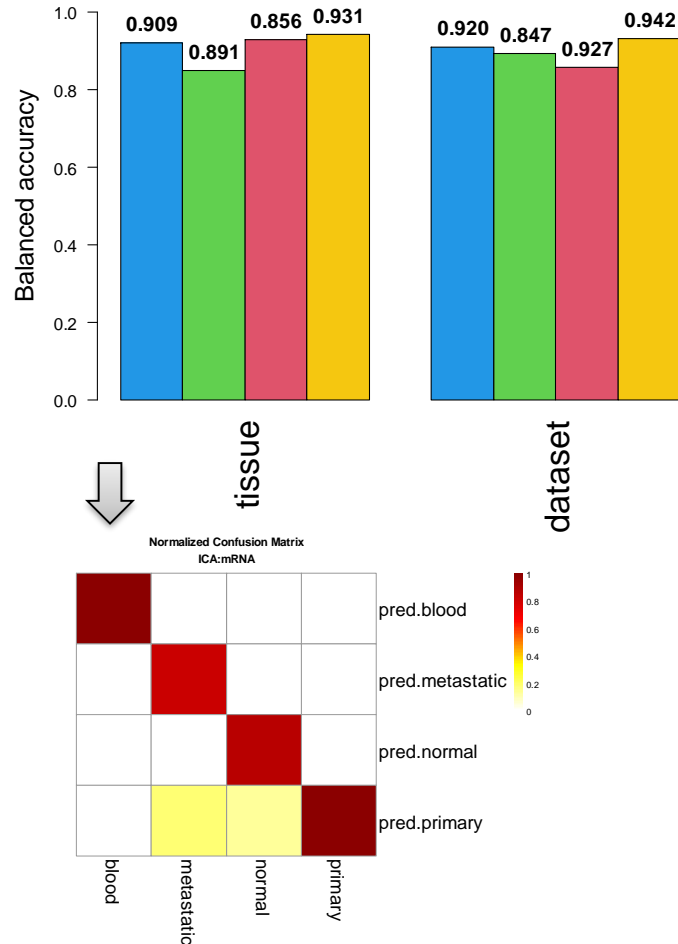
| Code | Study Name |
| --- | --- |
| ACC | Adrenocortical carcinoma |
| BLCA | Bladder urothelial carcinoma |
| BRCA | Breast invasive carcinoma |
| CESC | Cervical sq. cell carcinoma and endocervical adenocarcinoma |
| CHOL | Cholangiocarcinoma |
| COAD | Colon adenocarcinoma |
| DLBC | Lymphoid neoplasm diffuse large b-cell lymphoma |
| ESCA | Esophageal carcinoma |
| GBM | Glioblastoma multiforme |
| HNSC | Head and neck squamous cell carcinoma |
| KICH | 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 |
| THYM | Thymoma |
| UCEC | Uterine corpus endometrial carcinoma |
| UCS | Uterine carcinosarcoma |
| UVM | Uveal melanoma |

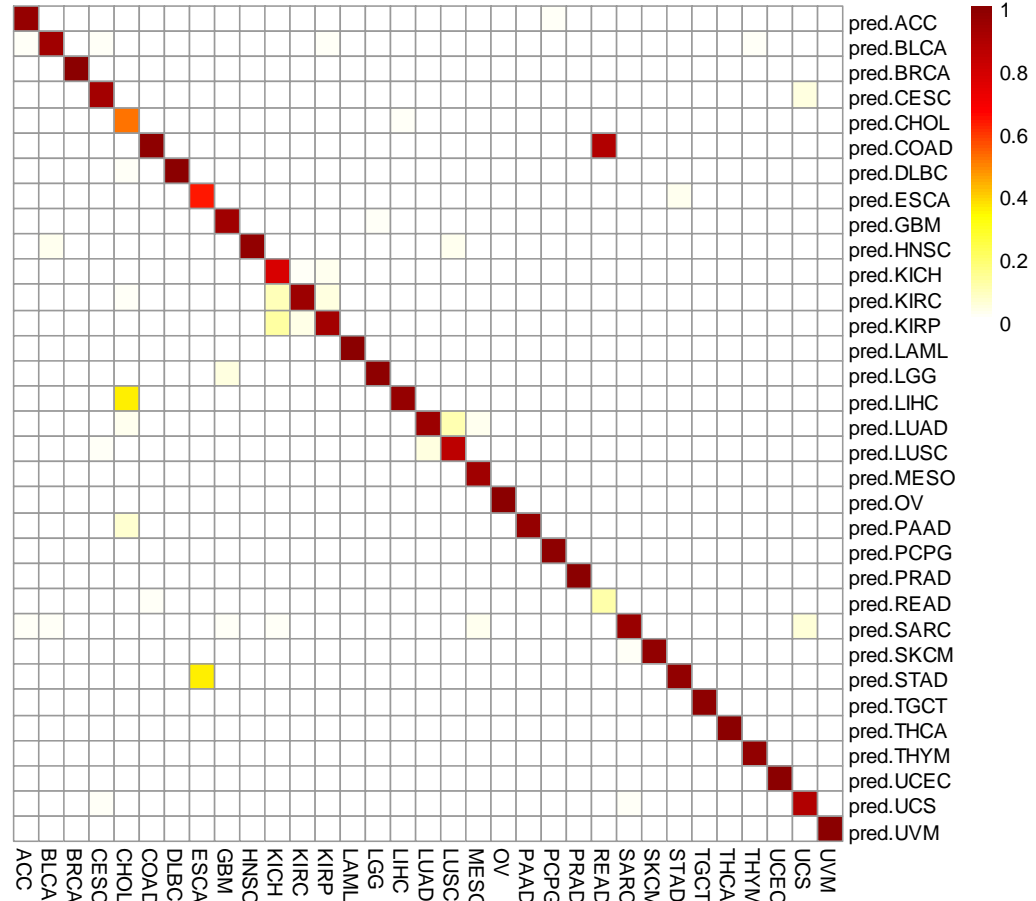# Pan-cancer: ICA-based Data Integration
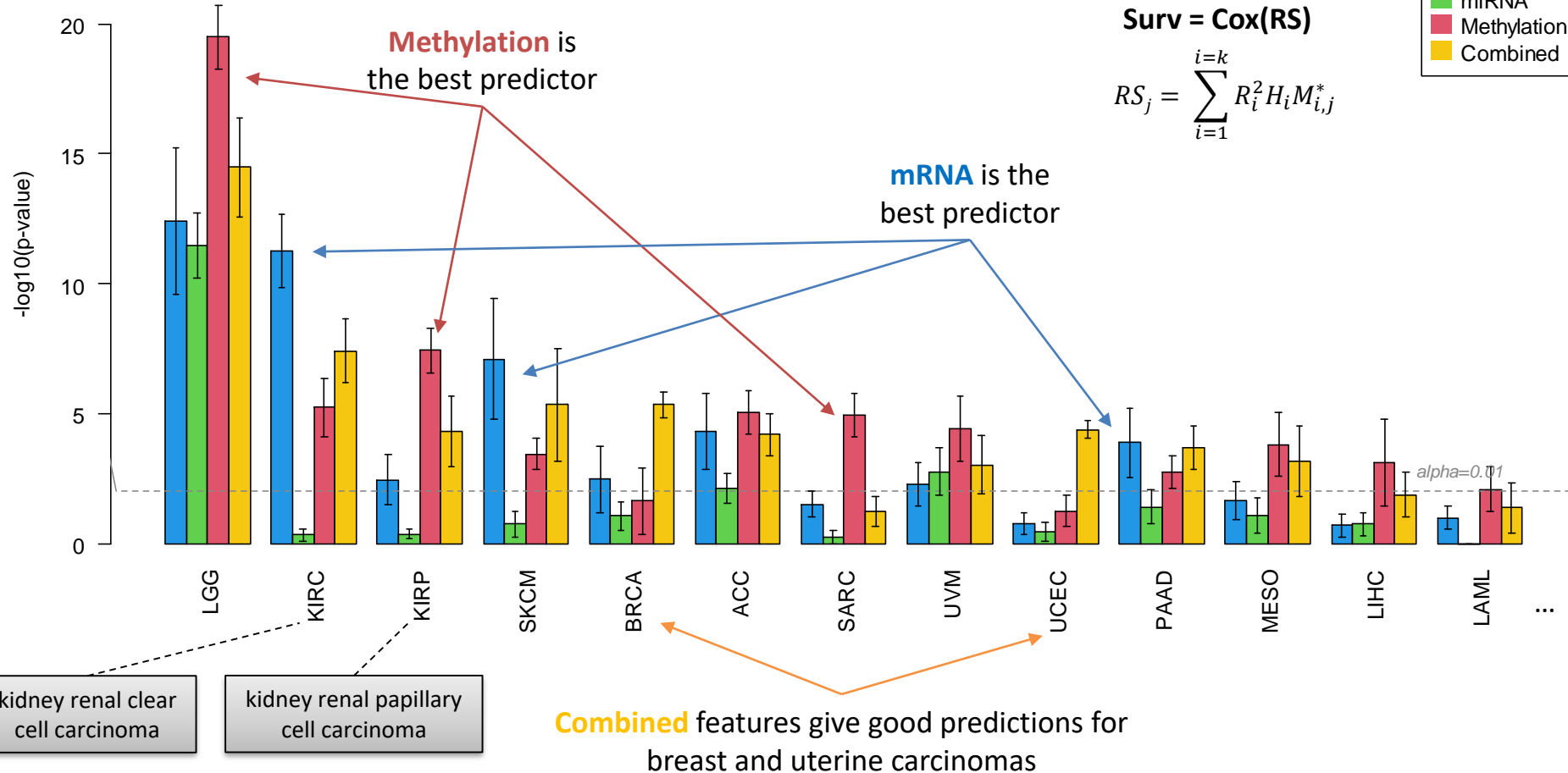
# Pan-cancer: Classification

# Pan-cancer: Prognosis

## Prediction of survival (same cohort, cross-validation)

Legend: mRNA, miRNA, Methylation, Combined

**Methylation** is the best predictor

**mRNA** is the best predictor

**Surv = Cox(RS)**

$$RS_j = \sum_{i=1}^{i=k} R_i^2 H_i M_{i,j}^*$$

y-axis: -log10(p-value), values 0, 5, 10, 15, 20

alpha=0.01

x-axis categories: LGG, KIRC, KIRP, SKCM, BRCA, ACC, SARC, UVM, UCEC, PAAD, MESO, LIHC, LAML, ...

kidney renal clear cell carcinoma

kidney renal papillary cell carcinoma

**Combined** features give good predictions for breast and uterine carcinomas

LUXEMBOURG INSTITUTE OF HEALTH

- ICA-based deconvolution:
  - ➤ 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
  - ➤ Reduce dimensionality

- Was validated:
  - ➤ Using acceptable computational methods (cross-validation)
  - ➤ On cell lines
  - ➤ Independent cohorts of patients

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

**(a) Omics Data Deconvolution** — done

**(c) Integration** — WP2
image features are used to predict component weights

**(b) Image Feature Extraction** — WP1

**(d) Validation and Control** — WP3-4

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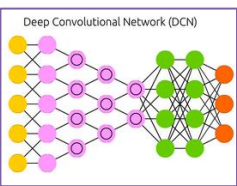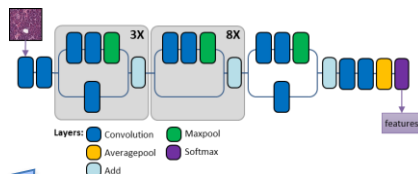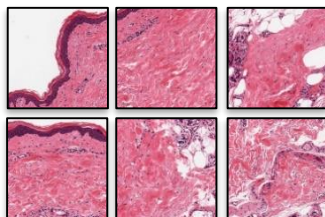
**(d) A thorough validation of the results** include (i) validation of an external pancreatic cancer cohort (DKFZ) and collection and (ii) in-depth analysis of in-house (LNS) samples of glioma patients. The expertise of the Co-PI (pathologist) will be used to validated predictions and the PI and his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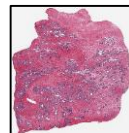
# Preliminary Results

# Tile-level Feature Extraction



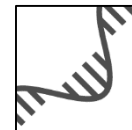**Examples of tiles classified with top certainty and co-localized with class medoids**



*Xception*, after parameter fine-tuning on organ classification task, transform each tile to ~150 non-zero features.
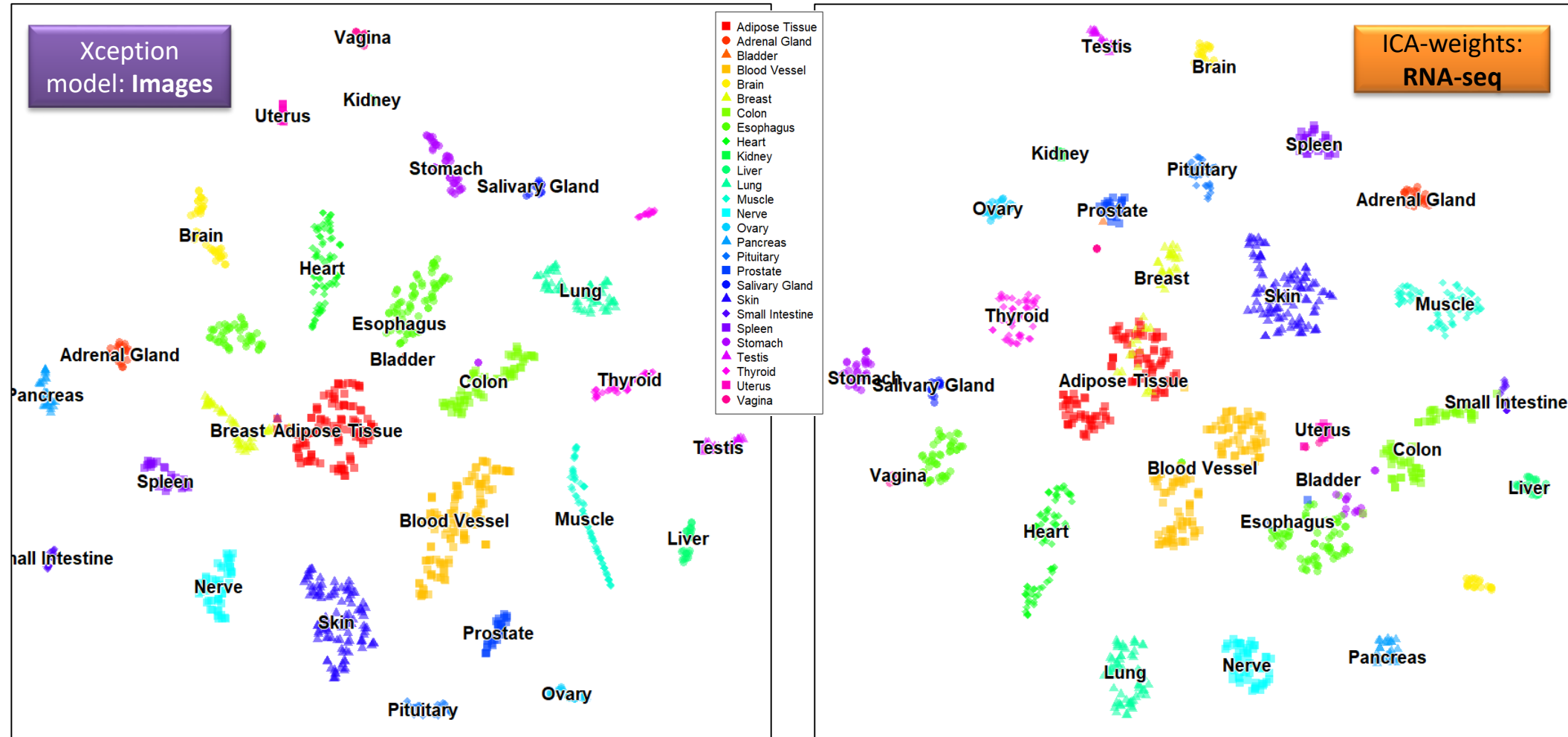
**Further analysis:**
These features were summarized to slide-level. Only 50% top-correlated tiles were preserved (can be further improved later...)
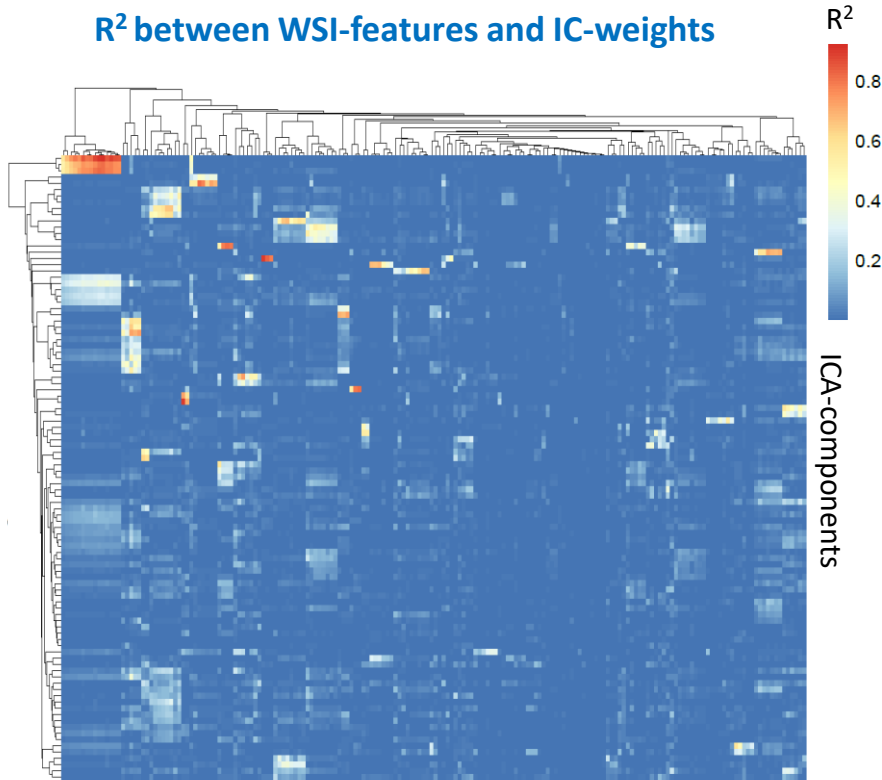
# Slide-level Analysis and ICA

## $R^2$ between WSI-features and IC-weights



WSI-features

ICA-components

## Predicting IC-weight



$R^2=0.9$

Brain

Pituitary

Predictions by RF on WSI-features

Weight of IC7 "synaptic transmission"

| GO:BP linked to IC7 | FDR |
|---|---|
| chemical synaptic transmission | 8e-28 |
| regulation of membrane potential | 8e-28 |
| behavior | 4e-22 |
| regulation of ion transport | 6e-22 |
| synaptic vesicle cycle | 3e-20 |
| cognition | 7e-20 |

## Predicting ICA-components

- 20% of the components were predicted with $R^2>0.9$
- 89% – with $R^2>0.5$



nature
COMMUNICATIONS

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

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

## Predicting genes

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

➢ 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)

➢ Combining molecular and his histopathological data may:

    ➢ Help pathologists faster and more accurate classify samples

    ➢ Improve accuracy of automatic data analysis

➢ Spatial transcriptomics, perhaps is our future ☺
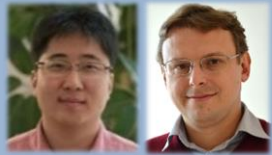
# Acknowledgements

## Bioinformatics Platform
### @ Data Integration and Analysis unit

R.Toth
V.Despotovic
L.Zhang
T.Koma
A.Muller

S-Y.Kim
P.Nazarov

LUXGEN

A.Aalto
Y.Zhang
B.Nosirov
T.Lukashiv

NORLUX @ DoCR

## Multiomics Data Science
### group @ Cancer Research

## Key internal collaborators

Simone
Niclou

Anna
Golebiewska

Michel
Mittelbronn

## Interns / students

Maryna
Chepeleva
(PhD student)

Aliaksandra
Kakoichankava
(PhD student)

Yibioa
Wang
(MSc)

Thomas
Eveno
(MSc)

Laurene
Picandet
(MSc)

## Key external collaborators

LSRU, Uni Luxembourg
Stephanie Kreis

Institute Curie, France
Andrei Zinovyev

DKFZ, Heidelberg
Jörg Hoheisel
Andrea Bauer
Nathalia Giese