

Motivation

Cancer dynamics is highly complex that data from a single layer of omics does not give the complete information to study a phenotypic characteristics of a cancer type.

While transcriptomic and proteomic abundance carry crucial information of cell states, studies have found that the two omics often do not agree with each other.

To better capture the biological context of interest, it is imperative to comparatively analyze them in the context of protein localization and biological pathway.

Although there are many data analysis tools for gene expression or protein quantification data, none of them allows researchers to compare different abundance data in the context of protein localization.

Methods: ALPACA

We present ALPACA (A Location-wise Proteome/transcriptome Abundance Comparative Analyzer) a visual data mining system that comparatively analyzes transcriptomic and proteomic abundance data of different cancers in location-wise and biological pathway-specific way.

Our system compartmentalizes the whole transcriptome and proteome abundance and visually presents the discrepancies of different cancers using subcellular locations and biological pathways as a filtering and sorting mechanisms.

Such filtering and sorting adds biological context of interest to the data analysis to aid the identification of potential biomarkers. ALPACA enables researchers explore the vast search space of numerous cancer types and pathway combinations with little effort and time.

Data Sets

Protein Localization Data

COMPARTMENTS, UniProt, Human Protein Atlas, LocTree3

34,172 unified protein-location labels

Proteomic / transcriptomic Abundance Data

CPTAC API

Transcriptome / Proteome for 11 Cancer Types/Subtypes

Pathway Data

KEGG pathway database

All Pathway Categories
All Constituent Genes for each Pathway

ALPACA Data Mining Process

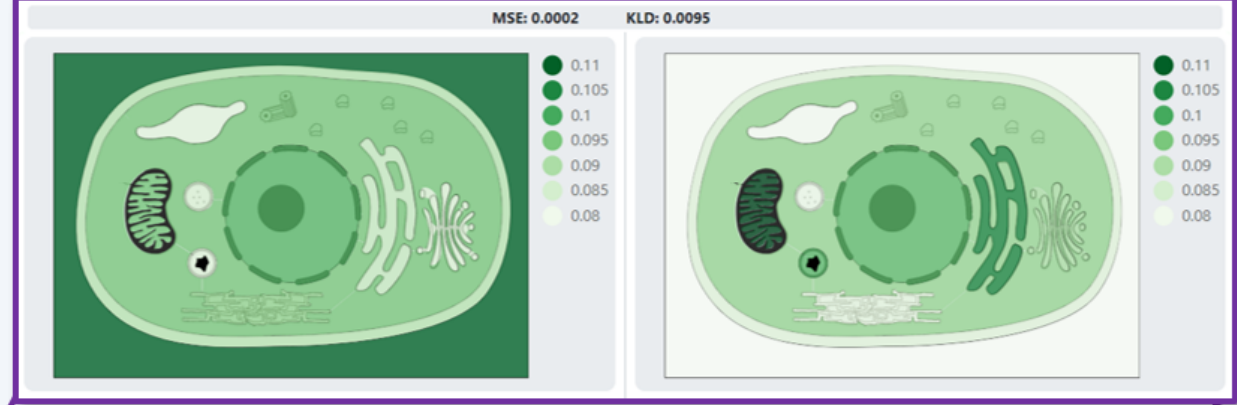
1. Select Omics & Cancer types

Choose Comparison Type: ☒ Proteomics
☐ Transcriptomics
☐ Combined(Trans vs Prot)

Choose Data 1:

Choose Data 2:

3. Location-specific Heatmaps



2. Select function(s) (pathway)

Sort by: ☒ Category ☐ Name ☐ OKLD ☐ OMSE

Select	KEGG pathway	KLD	MSE
<input checked="" type="checkbox"/>	0. Summary		
<input checked="" type="checkbox"/>	0.0 Summary		
<input type="checkbox"/>	Summary	0.0095	0.0002
<input checked="" type="checkbox"/>	1. Metabolism		
<input checked="" type="checkbox"/>	1.0 Global and overview maps		
<input type="checkbox"/>	Biosynthesis of amino acids	0.1355	0.0040
<input type="checkbox"/>	Biosynthesis of cofactors	0.1226	0.0032
<input type="checkbox"/>	2-Oxocarboxylic acid metabolism	0.0511	0.0014
<input type="checkbox"/>	Carbon metabolism	0.0460	0.0010
<input type="checkbox"/>	Fatty acid metabolism	0.0248	0.0005
<input type="checkbox"/>	Metabolic pathways	0.0127	0.0002

4. Identify Potential Key Proteins

symbol	location	kegg	exprs
A1BG	Extracellular space	Summary	2.2901
A2M	Extracellular space	Complement and coagulation cascades	2.6987
A2ML1	Extracellular space	Summary	2.1861
AAAS	Cytoskeleton	RNA transport	1.8529
AAAS	Cytosol	RNA transport	1.8529
AAAS	Nucleus	RNA transport	1.8529
AACS	Cytosol	Summary	0.6476
AADAT	Cytosol	Summary	1.4420

$$Q_p = \sum_{s \in S} \frac{\exp(q_{p,s})}{n}$$

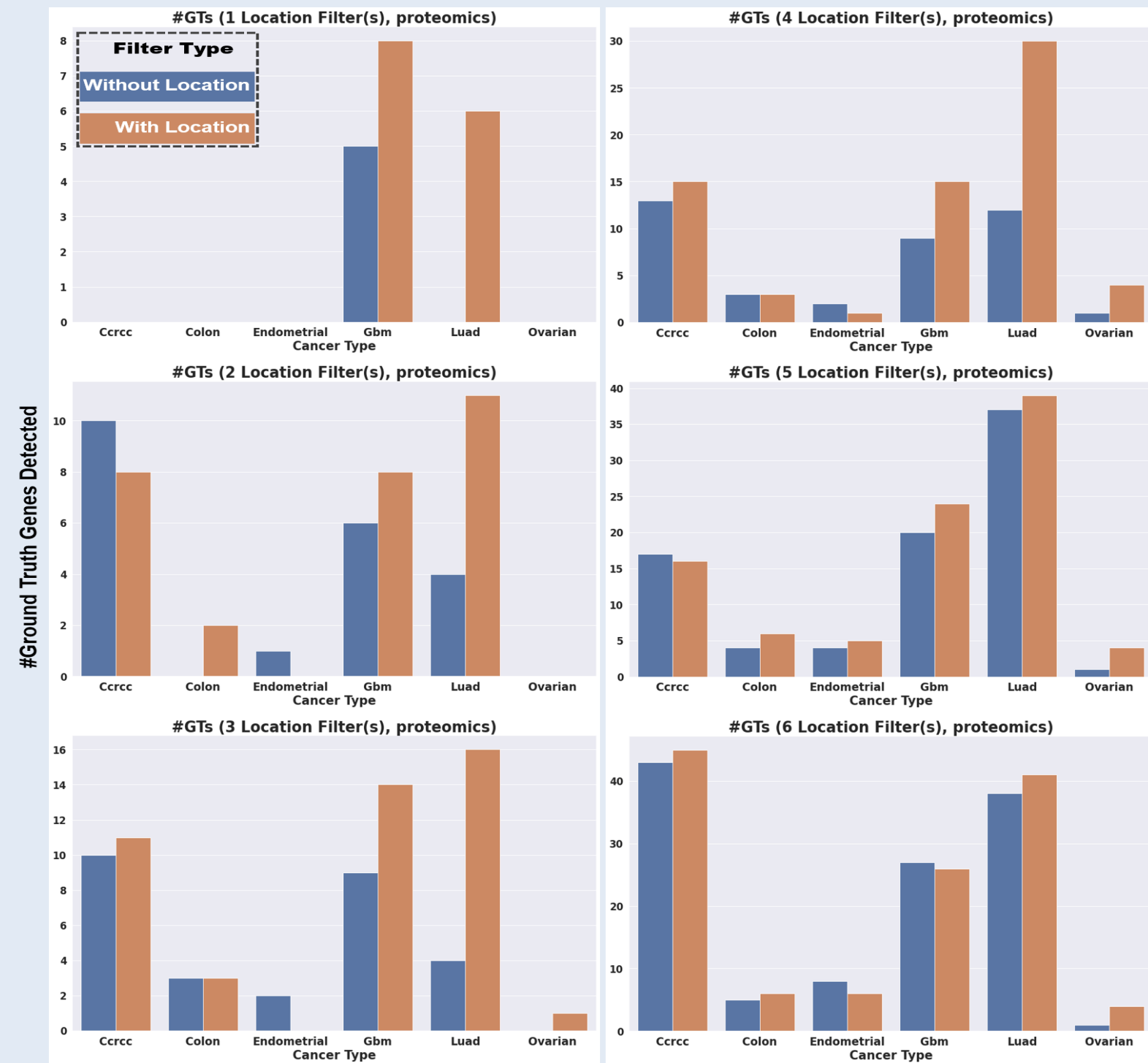
$$\mathbb{E}[Q_{p,\ell}] = \frac{\sum_p Q_p \cdot \mathbf{1}(p \in \ell)}{|\{p \mid p \in \ell\}|}$$

$$D(L)^\ell = \frac{\exp(\mathbb{E}[Q_{p,\ell}])}{\sum_{\ell} \exp(\mathbb{E}[Q_{p,\ell}])}$$

Results

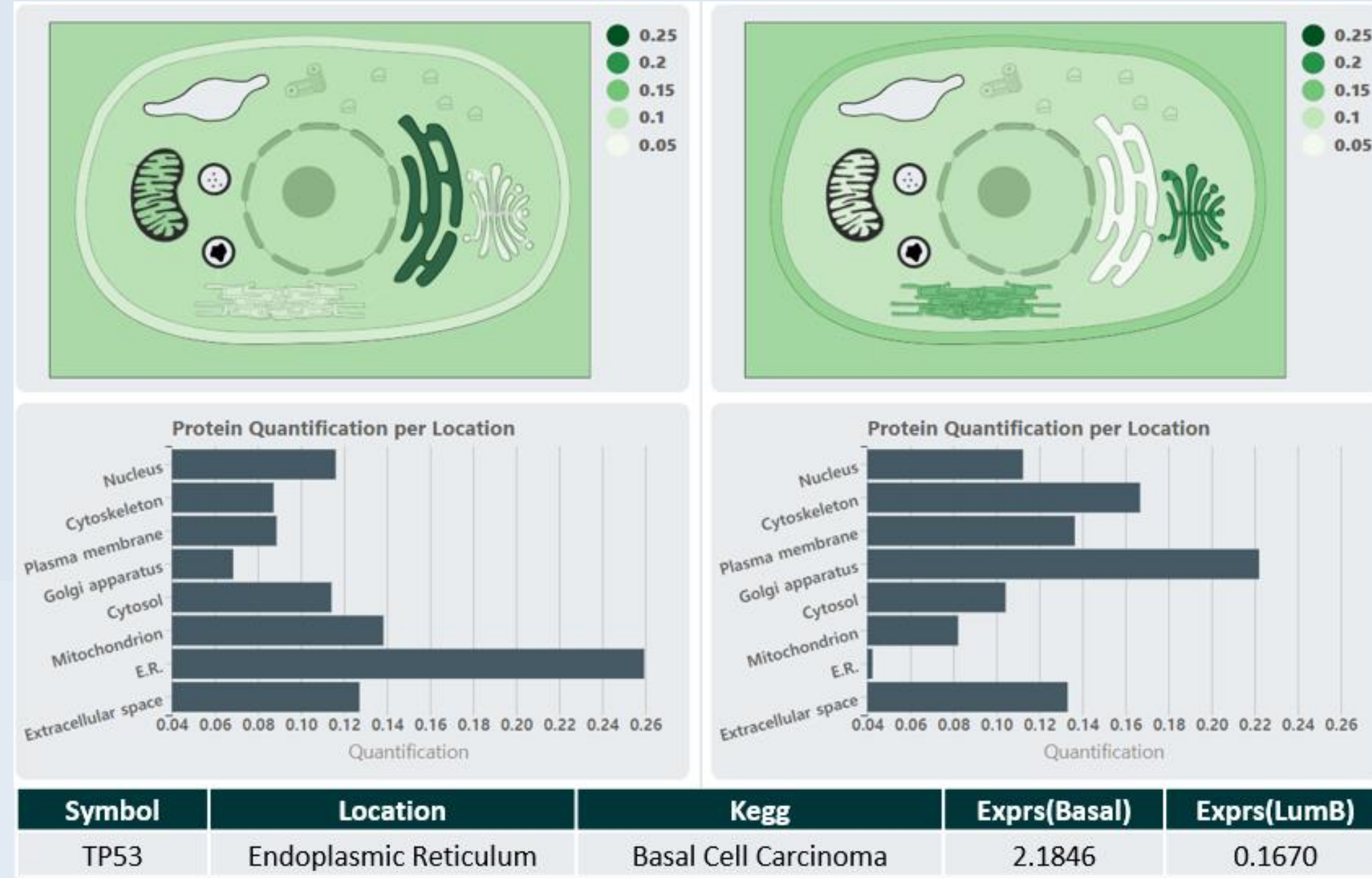
Power of Location Information

We demonstrate the utility of our system through cases studies of various scenarios. Guided by the quantitative difference metrics that ALPACA provides, we were able to efficiently search through different combinations of cancers and pathways to narrow down on the potential key proteins that can help explain the manifestation of a phenotypic difference between two cancers.



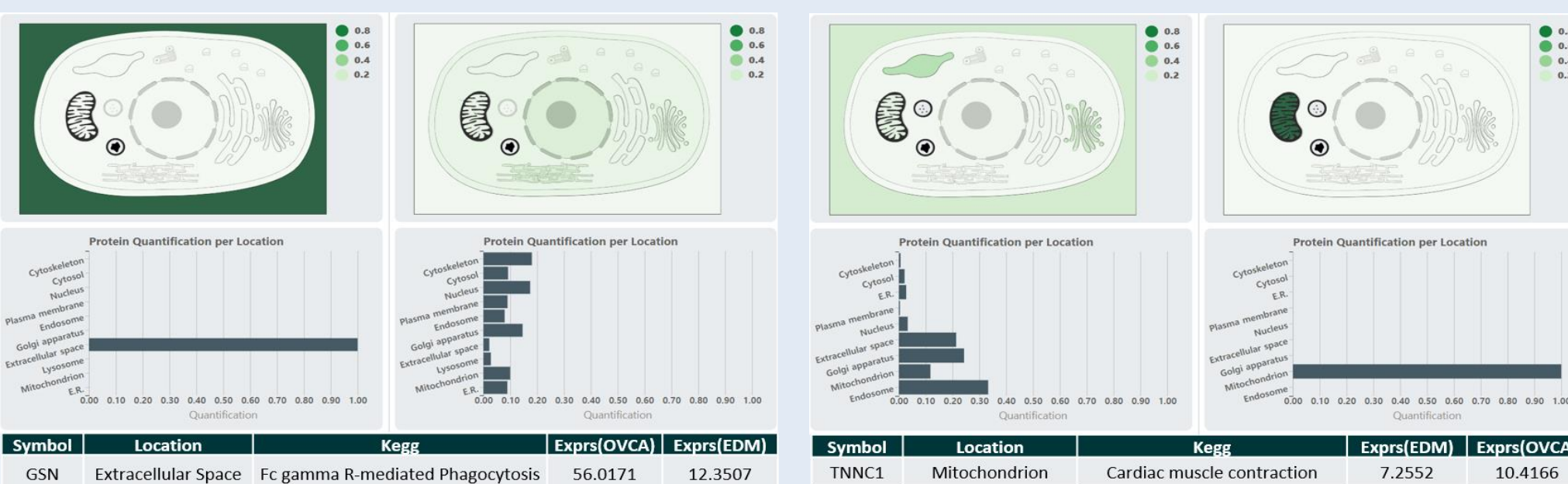
Data Mining Case Study

1. Rediscovery of Biomarkers



Largest discrepancy in E.R.: Overexpression of TP53 → TP53 implies its mutations. Known as a sign of poor prognosis. Supports the worst prognosis of Basal Subtype.

1. Data Mining Guided by Difference Metrics



Pathways with the largest location enrichment discrepancies according to the difference metrics reveal known biomarkers that may explain the phenotypic characteristics of Gynaecological cancers. GSN: known for poor prognosis in OVCA. ALPACA shows that GSN is indeed expressed the highest in OVCA compared to other Gynaecological cancers.