

Ruprecht-Karls-University Heidelberg  
Faculty for Life Sciences  
Molecular Biotechnology

# Thyroid cancer: Comparison of linear model and neuronal network (xxx) 3-Sätze-Zusammenfassung sfsf

Data Science Project SoSe 2022

Autoren Anna Lange, David Matuschek, Jakob Then, Maren Schneider  
Abgabetermin 20.07.2022

# Abstract

In the recent years bioinformatic methods became a tool of utmost importance in medical research. To define specific genes and pathways in different cancer types or histological types pan-cancer analysis are done. A focused analysis is done to specify different subcategories within a certain cancer type and to identify targets for targeted therapy. The main methods in identifying up- or down-regulated pathways are GSEA and GSVA. GSVA of TCGA expression data reveals four clusters of cancer types, which are defined by different histological types like glioblastoma and adenocarcinoma. The histological types therefore seems to correlate with a specific set of pathways being especially enriched in certain cancer types. Furthermore, a GSVA of Thyroid cancer expression data shows that thyroid carcinogenesis is associated with the up-regulation of proliferative signalling pathways like the hedgehog pathway and alpha6beta4 integrin signaling pathway and associated pathways such as IL-36 signaling. It also showed the down-regulation of a pathway that is associated with an increased MAP-kinase activity. It is based on those proliferative signalling pathways that three subclusters form inside of the THCA patients from the pan-cancer data. One THCA subtype that could be linked to the follicular histological subtype is defined by increased mTOR and MAPK activity, while having low alpha6beta4 activity. In contrast another THCA subtype is defined by a low mTOR and MAPK activity, but a high alpha6beta4 activity. The third THCA subtype is linked to enhanced activity of both of these proliferative signalling pathways. These results promise better results in treatment, as a more precise diagnosis of the distinct THCA subtype is possible. To improve the understanding of THCA and thereby hopefully improve patients prognosis, this project focuses on finding genes that have a significantly different expression in THCA compared to other cancers and especially to normal tissue.

# Abbreviations

xxxxxx

# Contents

Abstract	2
Abbreviations	3
1 Introduction	5
1.1 Biological background . . . . .	5
1.2 Computational tools . . . . .	6
1.3 Goals for the analysis . . . . .	8
2 Materials and Methods	9
2.1 Preprocessing of expression data . . . . .	9
2.2 Methods for descriptive analysis . . . . .	10
2.3 Dimension reduction and pathway enrichment analysis . . . . .	11
2.4 Regression analysis . . . . .	11
2.5 Environment . . . . .	12
3 Results	13
3.1 Preprocessing . . . . .	13
3.2 Descriptive analysis . . . . .	13
3.3 Pan cancer analysis . . . . .	15
3.4 Focused analysis . . . . .	19
3.5 Regression analysis of THCA pathway activity . . . . .	23
4 Discussion	25
4.1 Conclusion . . . . .	27
5 Outlook	28
6 References	29
7 Appendix	32

# 1 Introduction

## 1.1 Biological background

In 2019 230,000 cancer deaths were documented in Germany<sup>1</sup>. To detect and fight these tumors, the development of new treatment and detection methods is essential. A crucial step in this direction was the characterization of the Hallmarks of Cancer - properties present in every tumor. Among others those are: resisting cell death, inducing angiogenesis, enabling replicative immortality, activating invasion and evading growth suppressors (Hanahan and Weinberg, 2011).

However, the different histological types of tumors are equally important when characterizing cancers. Carcinoma, which can be further subcategorized into adenocarcinoma, squamous cell carcinoma, transitional cell carcinoma, derive from epithelial cells. Melanoma are skin tumors, sarcoma derive from mesenchymal tissue, glioblastoma from cells in the central nervous system and leukemia affect bloodcells (Alberts and Walter, 2015).

Here we focus on Thyroid carcinoma (THCA) as its incidence increased dramatically over the past few years (Cabanillas *et al.*, 2016). It arises in the thyroid gland which main function is synthesizing hormones and regulating body temperate and metabolism (Tsibulnikov *et al.*, 2020). THCA can occur in two different types, differentiated and undifferentiated thyroid cancer. Those two types again have histological subtypes. Papillary thyroid cancer (PTC), the most common THCA, follicular thyroid cancer (FTC) and a tall cell variant (TCV) are subtypes of differentiated thyroid cancer (DTC). Medullary and anaplastic thyroid cancer are subtypes of undifferentiated thyroid cancer (UTC). The prevalence of DTCs is clearly higher than of UTCs (Prete *et al.*, 2020). Among all THCAs DTCs, PTCs have the best clinical prognosis (Lin, 2007), while TCV cancers have the worst clinical outcome (Coca-Pelaz *et al.*, 2020). Therefore, the characterization of the differnt gene expression for each histological subype would be important for more specific therapy options. Even though, all thyroid cancers are treated with thyroidectomy and radioactive iodine, the additional therapy differs for each histological type (Kant *et al.*, 2020). For example, only some THCAs subtypes experiance integrin alpha6beta4 driven

---

<sup>1</sup><https://www.krebsinformationsdienst.de/tumorarten/grundlagen/krebsstatistiken.php>

carcinogenesis which might provide a viable target for therapy. Integrins are cellular adhesion molecules, that bind to laminin in the extracellular matrix (Liberzon *et al.*, 2015a). Together with other proteins they form hemidesmosomes. Thereby, integrins are essential for the integrity between cells. An important step in the development of malignant cancer is the invasion into healthy tissue. Thus, the detachment of the extracellular matrix from the surrounding cells is essential and alterations of integrin are very common in cancer cells (Rabinovitz and Mercurio, 1996).

## 1.2 Computational tools

To analyse how the activity of a gene set differs between two sets of gene expression data, a Gene Set Enrichment Analysis (GSEA) is performed. For this, the genes in the expression data have to be ranked decreasingly by a certain metric. Such metrics can include the log<sub>2</sub> fold change between the sample expression data and a reference set or the associated p-values for each gene. After ranking, a cumulative sum of all expression values in the ranked sample is computed. If a gene is present in the gene set to be analysed the expression value of that gene is added to the running sum. However, if the current gene does not lie in the gene set the value is subtracted. The extremum of this running sum is termed the enrichment score of the gene set. It is positive if the gene set is overexpressed in the sample compared to the reference data and negative vice versa. (Reimand *et al.*, 2019)

The Gene Set Variation Analysis (GSVA) is performed with the same intention as the GSEA - to analyse the gene set activities in gene expression data. However, no reference data is required to successfully perform GSVA. There are various approaches to GSVA, one of them is performed by (Hänzelmann *et al.*, 2013a) by following five steps. First, the cumulative density distribution of a gene over all samples is estimated. Then the expression statistic of a gene in a sample based on the cumulative density distribution is calculated to bring all of the expression values to the same level. The third step is to rank the genes based on the expression statistic and to normalize the ranks with z-transformation. Lastly, the enrichment score is computed based on the obtained ranked list by calculating the Kolmogorov-Smirnov-like rank statistic for each gene set. (Hänzelmann *et al.*, 2013a)

A Principal component analysis (PCA) xxx QUELLE is used to alter the coordinates of a given dataset to its eigenvectors. This matrix rotation results in a new set of basis vectors called principal components (PCs) - the eigenvectors - that are orthogonal and show little correlation. Sorting the PCs by their associated eigenvalue, the PCs explaining the most variance can easily be identified, as they have the highest eigenvalue. By displaying the data set in a coordinate

system span by the  $n$  most variant PCs, the dimensionality of the dataset is reduced to  $\mathbb{R}^n$  with the lowest loss in variance.

The Uniform manifold approximation and projection for dimension reduction (UMAP) is a method to reduce the dimension of a multidimensional data set. Compared to PCA, UMAP preserves the global structure of the data better and is much faster than other comparable techniques like t-SNE (Maaten and Hinton, 2008). The algorithm starts by setting up a high-dimensional graph representation of the data. From each data point, a radius is extended and when two radii come into contact the points are connected in the graph. The radius is chosen individually for each point based on the distance to the nearest neighbor. The algorithm goes on until  $k$  points are connected or  $n$  iterations are reached. The resulting clustered high-dimensional graph is then optimized for a visualization in low-dimensions. A disadvantage of UMAP is that although the overall structure is conserved, the distances between the individual points are not proportional to the real distance in the data set. This arises from the non-linear dimensional reduction. (Sharma *et al.*, 2021)

The Jaccard index is the intersection, divided by the union of two sets. Therefore, it can be used to identify the similarity of the sets.

Linear regression is a statistical model that uses measurable values to predict an outcome. For this purpose, a linear function serves as basis to build the linear regression equation (Lunt, 2013). The coefficients for each variable are estimated by their correlation and slope with the predicted parameter. Lastly, all coefficients as well as the intersect are optimized for the data set with a least sum of squares method.

As an alternative to the linear regression a neuronal network can be used. In general, a deep learning network consists of an input layer, multiple hidden layers, and an output layer consisting of various neurons (Riedmiller). The input layer contains as much neurons as input numbers are given for each sample. The output is for a regression analysis a singular neuron. The number of neurons in each hidden layer and the number of hidden layers vary and must be tested to give best results. The activation of each neuron can be described as a linear composition of all the inputs  $x_i$  from the previous layer associated with a weight  $w_i$  and a bias:

$$Activation = \sum_{i=1}^n w_i x_i + bias$$

To obtain neuron activations in the range of 0 and 1, a min/max-scaling is performed on the input data. The “learning effect” of the network is achieved by optimizing the randomly chosen

weights and biases via gradient decent. To do so, for each training iteration the error of the network is computed by a cost function:

$$Costfunction = \frac{1}{2m} \sum_{i=1}^m (x - y)^2$$

$m$  is the number of samples,  $y$  the output and  $x$  the expected output.

Next, the cost function value must be reduced. Therefore, its gradient is computed, and all weights and biases are adjusted accordingly in a process called backpropagation. In resilient backpropagation, only the sign of the gradient is used, to avoid harmful effects of its magnitude. For the next samples those steps are repeated to reach the minimum of the cost function. A drawback of this method is that gradient decent only identifies local minima of the cost function. To find a global minimum the training has to be repeated with various initial weights and biases. After such a minimum is identified to network performs optimally for the data set.

### 1.3 Goals for the analysis

Using pan-cancer expression data we are going to identify clusters between cancer types regarding their expression profiles as well es hallmark and metabolic pathway activities. Then, our focus will shift to THCA, were we are going to identify subclusters in gene expression linking them to histological types. Additionally, we are going to analyse pathways that alter significantly between THCA and homeostatic thyroid tissue and predict their activities with linear and neuronal network regression.



## 2 Materials and Methods

In the course of this project two separate analysis are performed: a pan-cancer analysis focusing on differences between cancer types and a focused analysis investigating THCA.

For the analysis four data sets were provided. For pan-cancer analysis a gene expression data frame with normalized and log2 transformed bulk RNA-seq expression data for 60,489 genes in 9741 patients with 33 different forms of cancer was used. The data was derived from The Cancer Genome Atlas (TCGA). Complementing the TCGA expression data is an annotation data frame with 37 clinical annotations regarding tumor type, tumor stage, gender, age, etc. for all patients.

The third piece of data is a list containing five lists for the focused analysis, one list for each tumor type (BRCA, KIRC, LUAD, PRAD, THCA). For our focused analysis, only the THCA data were used. The THCA list consists of three data frames: The first two contain normalized and log2 transformed bulk RNA-seq expression data for 19,624 genes in 59 THCA patients for carcinogenic and homeostatic tissue. The third data frame complements the data with the respective clinical annotations.

The last object contains 46 pathways associated with the hallmarks of cancer in form of a list of string vectors.

To perform enrichment analysis later on, 6366 canonical pathways were selected from the Molecular Signatures Database (MSigDB)(Liberzon *et al.*, 2015b) with the `msigdb::msigdb()` function. As not to introduce a bias during enrichment analysis, the similarity of MSigDB pathways among themselves as well as with the hallmark pathways was computed with the Jaccard index. Pathways with a Jaccard index greater than the  $1\sigma$  range were discarded.

### 2.1 Preprocessing of expression data

All expression data were checked for missing values with the `na.omit()` function. Subsequently, low variance filtering was performed for TCGA and THCA tumor expression data. The variances

of expression were computed for every gene across all samples and then, genes with variances below a threshold were discarded to reduce dimensionality.

Next, biotype filtering was performed for pan-cancer and THCA expression data to reduce dimensionality further. Only genes sharing biotypes with the hallmark pathways were kept for the following analysis. The biotypes of the genes were retrieved using the `biomart::getBM()` function from the biomaRt package (Durinck *et al.*, 2009). To allow for an appropriate comparison within all pathways, only MSigDB pathways in which over 99% of their respective genes were present in the filtered expression data were selected as final pathways.

## 2.2 Methods for descriptive analysis

In a mean-variance plot the variance is plotted over the mean of expression values of single genes across all patients. Thus, the variance and mean were calculated for each gene in the THCA expression data. The final plot was created with the package `ggplot2` ??.

Jaccard index is a method to describe the similarity between two quantities. To compute it, the intersection of all gene EnsembleIDs from two compared pathways was divided by their union. We used this method to determine the degree in which pathways are similar to each other.

A volcano plot is used to identify genes displaying significantly different expression in cancerous versus homeostatic tissues. First, the log2 fold change (Log2FC) is calculated for each gene across all samples in the THCA expression data in the following way:

$$\log_2FC = \text{mean}(\text{normaltissue}) - \text{mean}(\text{tumortissue})$$

Next, a two-sided t-test was performed with the `t.test()` function to determine the significance of a difference in expression. To avoid the accumulation of type one errors, a Bonferroni correction was performed.  $n$  is the number of genes in the cleaned data set for focused analysis:

$$\alpha = \frac{0.025}{n}$$

In the volcano plot the  $-\log_{10}$  of the calculated p-values is plotted against the Log2FC. Genes with a lower p-value than the corrected significance level  $\alpha$  are significantly differently expressed. If the Log2FC is additionally positive, the genes are significantly overexpressed in tumor tissue, if the Log2FC is negative, the genes are significantly underexpressed in tumor tissue.

## 2.3 Dimension reduction and pathway enrichment analysis

The GSEA was used to identify enriched pathways in THCA tumor tissue. Here, GSEA was performed with the package “fgsea” (Korotkevich *et al.*, 2019). First, the expression values were ranked in decreasing order by Log2FC for every patient. Log2FC was chosen as the ranking metric as it is easy to compute and shows a high sensitivity. **xxx Quelle: Ranking metrics in gene set enrichment analysis: do they matter?** Secondly, using the ranked Log2FC vectors, the enrichment score of each pathway was calculated for each patient with the `fgseamultilevel()` function.

As no normal tissue reference data was provided for the TCGA expression data, pathway activities were computed via GSVA. The analysis was performed with the `gsva()` function from the “GSVA” package (Hänzelmann *et al.*, 2013b). To give a general overview over the differences in expression of THCA and homeostatic thyroid tissue GSVA, the THCA expression data were also analysed by GSVA. To do so, tumor and normal expression data were combined into a singular data frame of which enrichment scores were computed with `gsva()`. Then, the GSVA data was split again and the log2FC between the two matrices was computed and taken as pathway activity.

PCA was performed to provide an uncorrelated data set for the subsequent UMAP. For the TCGA GSVA pathway activity data the `prcomp()` function was used. To verify the results, PCA was performed on TCGA expression data, as well. In this case `Seurat::RunPCA()` from the Seurat package was used to minimize computation time (Hao *et al.*, 2021).

UMAP analysis was done on the principle components from previous PCA to identify and visualize clusters in TCGA GSVA and expression data. This was achieved with the `umap()` function from the package “umap” (Konopka, 2022) running on all PCs from TCGA GSVA and expression data. The computational effort is lower in UMAP than in PCA. UMAP works on uncorrelated features provided by the PCA.

## 2.4 Regression analysis

For THCA pathway activity regression analysis a highly variant and significantly altered pathway was selected. To prepare the data appropriately the THCA GSEA data set was divided into a training and test data set containing 44 and 15 samples respectively. A linear regression analysis was performed on the training data with the `glm()` function. To do so, the correlation of all pathways was computed and pathways with high correlations are omitted. Subsequently, the 10% of most variant pathways are selected as variables for the regression model. A second

model was introduced by computing the p-values of all coefficients and selecting only those pathways contributing significantly to the model were kept.

A neural network was implemented to predict the pathway activity using the `neuralnet()` function from the “neuralnet” package (Fritsch *et al.*, 2019). For identification of the best initial conditions, 25 different networks are generated, each with 2 hidden layers and different combinations of neurons per layer. For each combination the networked was trained on the min-max-scaled training data and the best network was determined by the lowest mean squared error (MSE) in the test data.

## 2.5 Environment

The R version 4.0.1 was used, the table of used packages is attached in the appendix (see table @ref(tab:packagesused)).

## 3 Results

### 3.1 Preprocessing

#### **Dimension reduction through low-variance and biotype filtering**

All the TCGA and THCA expression data were checked for NAs, which were subsequently deleted. Genes with a variance lower than 0.1 were removed to reduce dimensionality, as they contribute very little to the overall variance of the data set and are most likely house-keeping genes (xxx Quelle oder löschen). The low-variance filtering of the THCA data set was done in a similar way. Genes with a lower variance than 0.06 were deleted in the tumor tissue and the normal tissue data. To reduce dimensionality further, the biotype of the hallmark pathway genes was determined, which was almost exclusively protein coding. To match this, only protein coding pathways were kept in all expression data sets for further analysis. Doing so, the number of genes in the pan-cancer data set was reduced from 60,000 to approximately 19,000 genes and from approximately 20,000 genes to 15,000 genes in the THCA data.

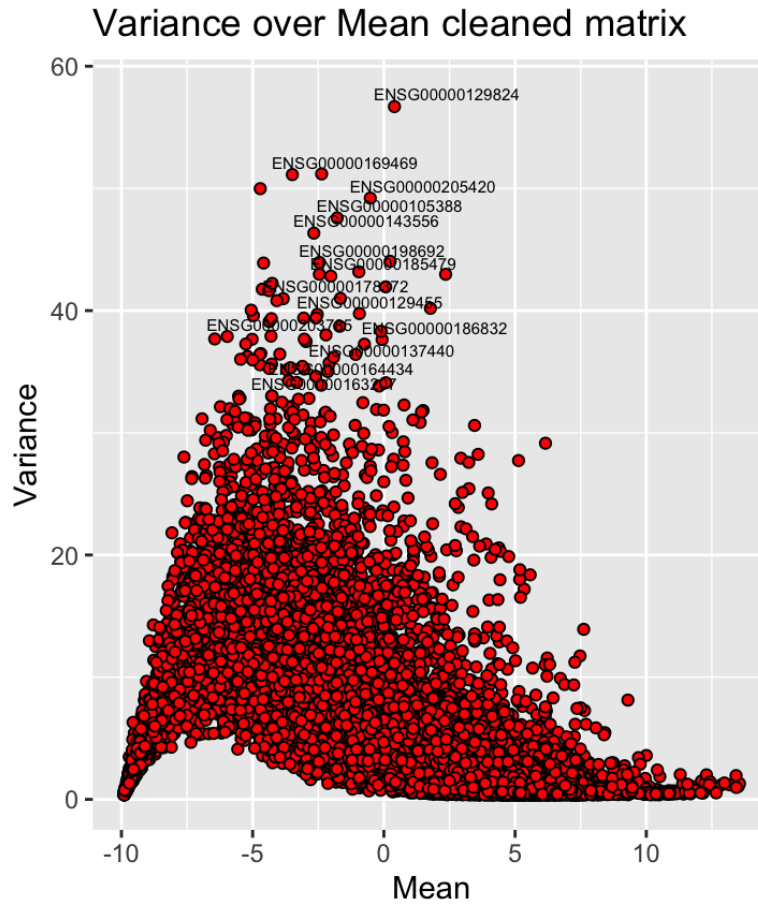
#### **MSigDB pathway filtering**

The pathways from the MSigDB database were first aligned with the genes in our expression data. Only pathways with a coverage of over 99% were kept. To test for similarity in the selected metabolic pathways compared to the hallmark pathways and the metabolic pathways themselves, the Jaccard index between all pathways was calculated. A few MSigDB pathways with a high Jaccard index were identified and subsequently deleted. Heatmaps, displaying the jaccard index of cleaned and uncleaned data can be seen in the appendix @ref(fig:Jaccarddirtea) @ref(fig:Jaccardcleaned). The number of MSigDB pathways could thus be reduced from 6366 to 657.

### 3.2 Descriptive analysis

**Mean-variance plot of TCGA expression data shows highly variant genes.**

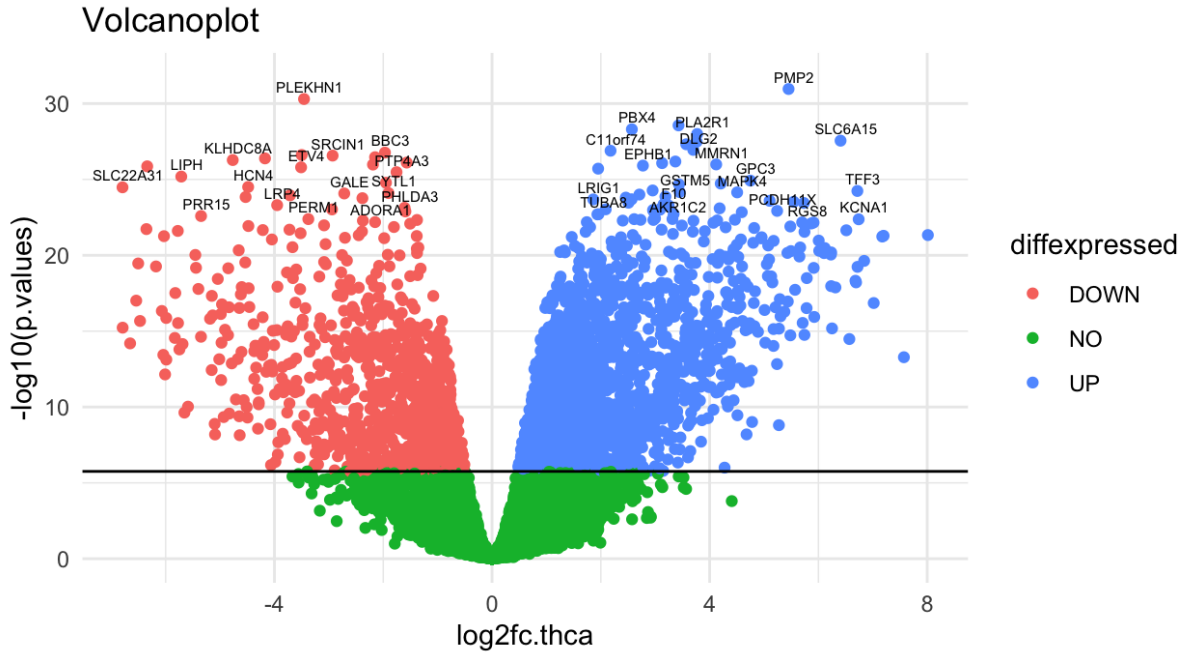
To determine the genes from the TCGA expression data with a high variance, the variance was plotted over the mean (Figure @ref(fig:showmeanvariance)). Additionally those genes with a variance higher than 33 were labeled with their EnsembleID. The distribution of genes in this plot shows that the highly variant genes are around a log2 mean expression level of 0. The plot also shows, that very few genes are at a low mean expression level or at a very high mean expression level. Most genes are expressed across all patients at a log2 mean expression level of approximately 0. With this plot we were able to determine which genes differ significantly in their expression level across all cancer patients.



**Figure 3.1:** Mean-variance plot of cleaned TCGA expression data. Y-axis shows variance of a gene expression, x-axis shows the log2 mean of a gene expression. Genes with variance greater 33 are labelled with their ENSEMBL-ID

### Significantly up- and down regulated genes in THCA obtained from volcano plots

To determine up- or down-regulated genes in THCA corresponding p-Values were computed with a Wilcoxon rank sum test. (Figure @ref(fig:showvolcanoplot)). The significance level was adjusted to 1.755e-06 with a Bonferroni adjustment.



**Figure 3.2:** Volcano plot of THCA expression data. Downregulated genes are colored red, upregulated genes blue. Not significantly altered genes are colored green. Most significantly altered genes are labelled with their gene symbol

### 3.3 Pan cancer analysis

#### GSVA of TCGA expression data reveals four clusters of cancer types.

A gene activity matrix was computed as described in section (methods). The obtained pathway activity matrix was visualised in a heatmap.

To find general clusters in the heatmap, the mean expression of each gene in each tumor type was generated and clustered hierarchically. Figure @ref(fig:meanexp) and @ref(fig:exp)

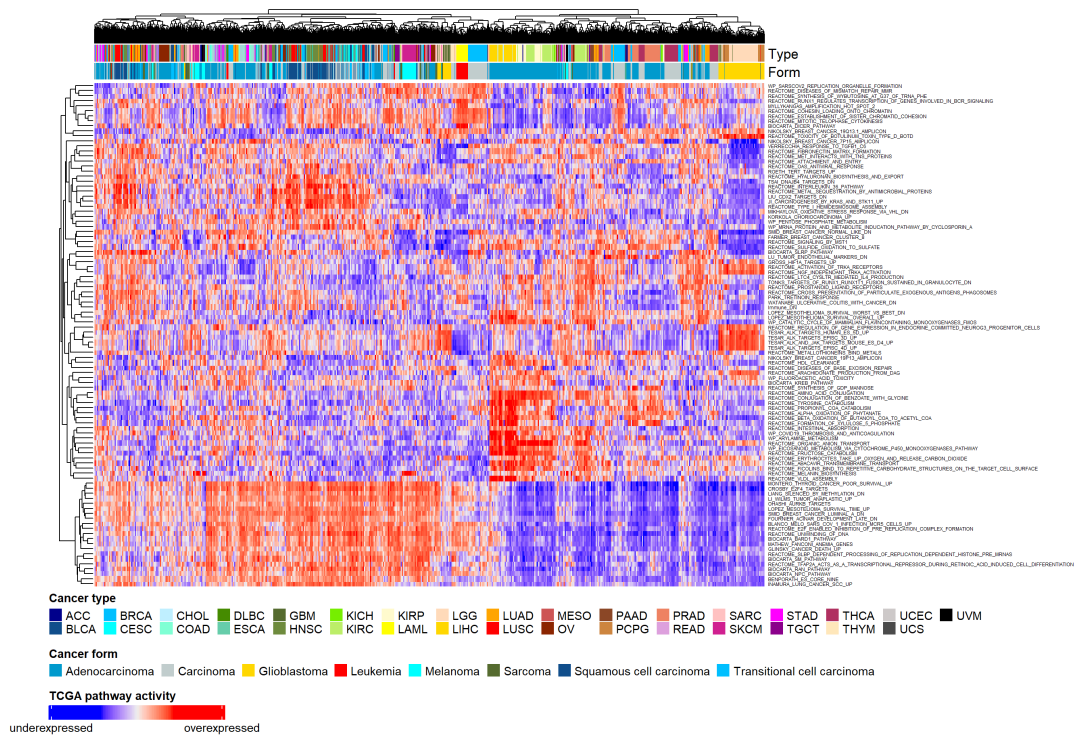
The tumor types were clustered based on their mean pathway activity and formed four clusters correlating with their histological type. The first cluster contains mainly adenocarcinomas, while the second one contains predominately glioblastomas. Leukemias are only found in the third cluster and the last cluster is enriched with sarcomas and carcinomas. Melanomas appear in the second and fourth cluster. Furthermore, three observations were made regarding specific information about pathway activity.

Pathways, which are important for nucleus import and export like Nasopharygeal carcinoma (NPC) and Ran shuttle pathways, as well as pathways for transcription regulaturs in embryonic





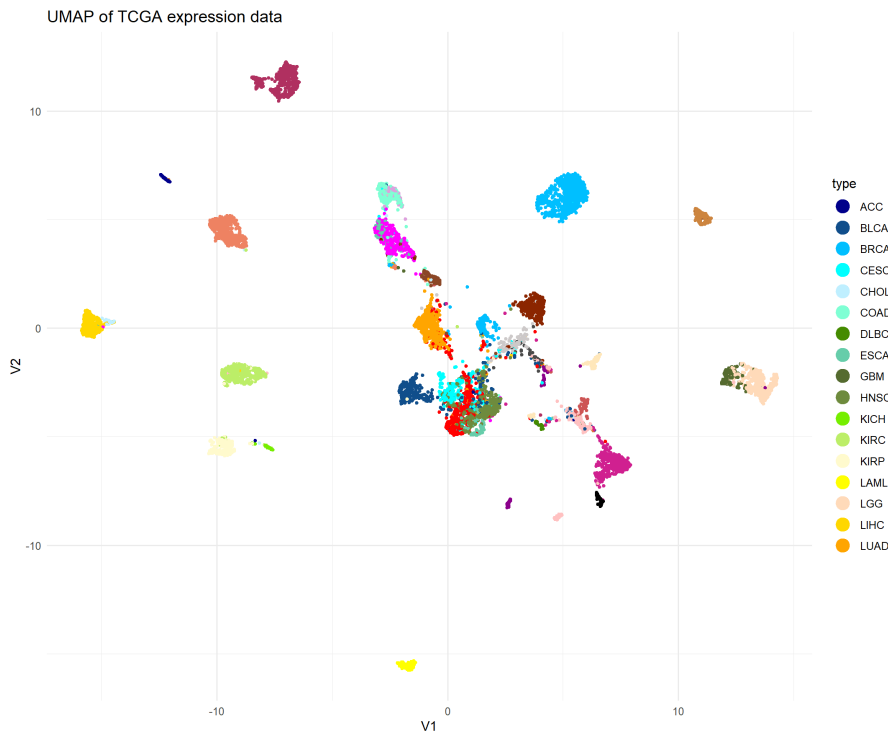
## RESULTS



**Figure 3.4:** Pathway activity of the 100 most variant pathways for each patient. Column and row clusters were obtained by complete hierarchical clustering. Pathway activities were computed via GSVA of pan-cancer expression data. For all pathway activities see figure (XXX in the appendix).

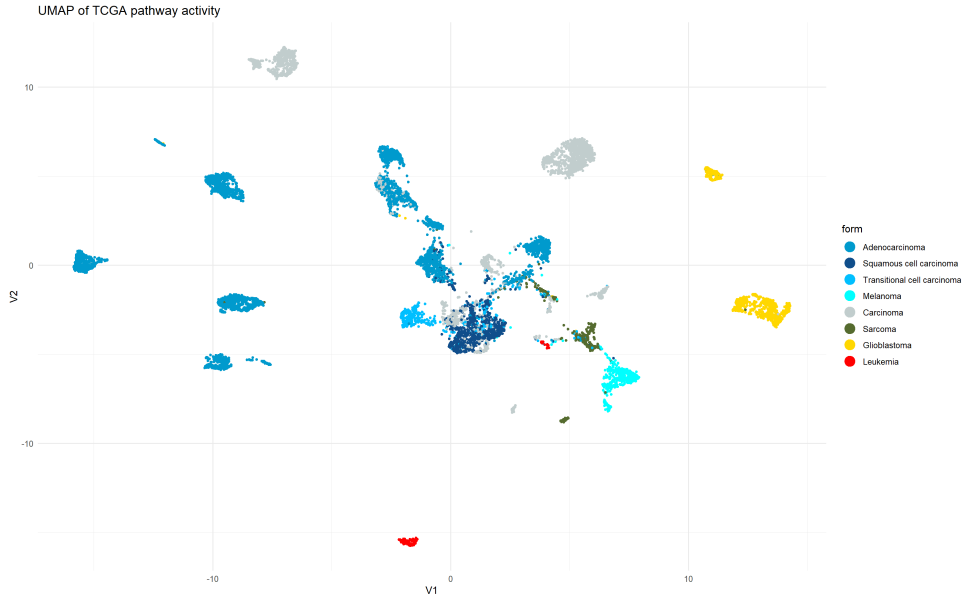
### Dimension reduction of GSVA pan-cancer data reveals clusters in pathway activity.

PCA was performed on GSVA pan-cancer for UMAP analysis. No apparent clustering was observed only in PCA data (compare Figure (ref?)(fig:PCAcancerform) @ref(fig:PCAcancertype) ). Subsequent UMAP analysis however, showed clear clusters for most cancer types. @ref(fig:UMAPPanType) @ref(fig:UMAPPanForm). This complements the results obtained from our heatmap and reassures, that the tumor types have characteristic pathway activities. However, some cancers cluster better with their histological type rather than tumor type. This was observed mainly for carcinomas like squamous cell carcinoma and transitional cell carcinoma, as well as sarcoma, lung adenocarcinoma and ovarian cancer. These are the same histological types that proofed difficult to cluster in the mean GSVA of TCGA expression. The UMAP confirmed the assumption, that the histological type of a tumor has a major impact on the patients gene expression profile.



**Figure 3.5:** UMAP of TCGA pathway activity, colored by tumor type

The same analysis was performed for gene expression activity instead of pathway activity to check for reliability of the results. Similar clusters were observed, which confirms our results. See @ref(fig:UMAPGenform), @ref(fig:UMAPGen) in the appendix.



**Figure 3.6:** UMAP of TCGA pathway activity, colored by histological type

### 3.4 Focused analysis

#### **GSVA on THCA expression data reveals pathways driving thyroid carcinogenesis.**

To grasp a general overview of the differences in pathway activity between THCA and homeostatic thyroid tissue, GSVA was performed for the THCA expression data. Then, changes in pathway activity were computed by log2 fold change and the respective p-values were computed by a Wilcoxon rank-sum test. The most significantly altered pathways were then characterized. @ref(fig:THCAvolcano) Most prominently among them were pathways linked to proliferative signaling such as upregulation of p53 inhibitory proteins and hedgehog pathway activating Gli proteins. Further, the alpha6beta4 integrin signaling pathway and associated pathways such as IL-36 signaling and Typ I hemidesmosome synthesis were significantly enhanced in THCA. Further, signaling through the EWSR1/FLI1-fusion protein was significantly upregulated in THCA. Lastly, THCAs showed downregulation of non-histone protein methylation.

#### **Pan-cancer data GSVA reveals three subtypes of THCA altering in proliferative signaling.**

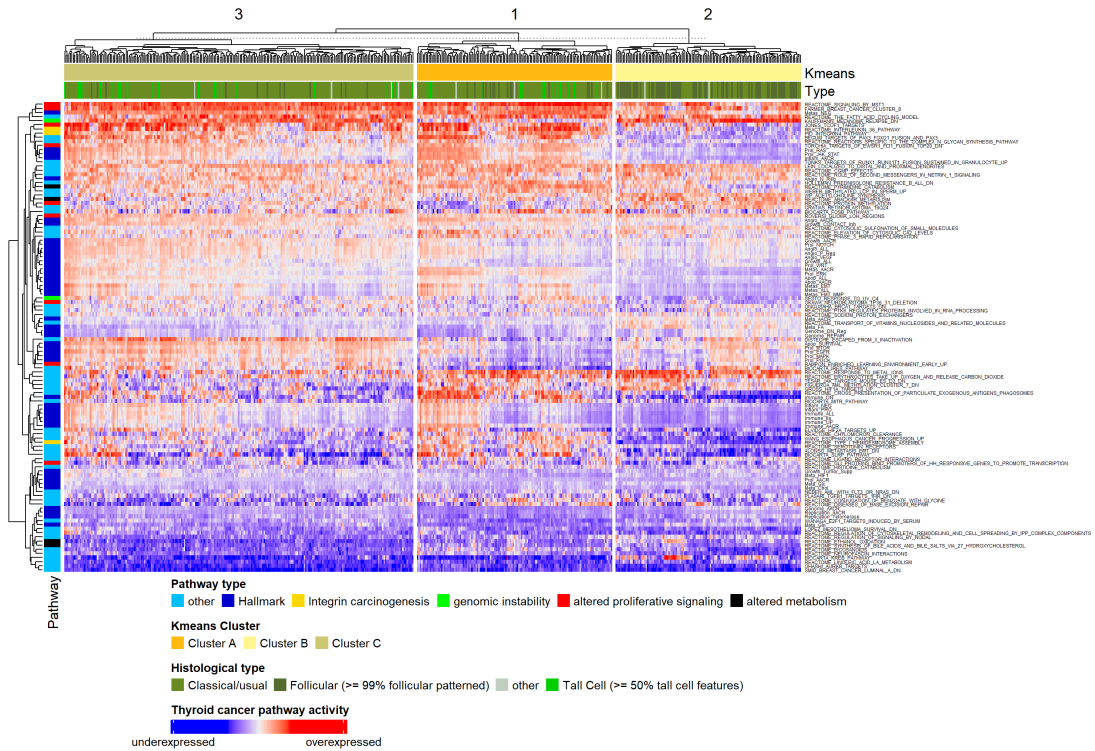
To investigate potential subtypes of THCA, the respective samples were taken from the pan-cancer GSVA data. The optimal number of clusters was determined by an elbow plot and subsequent K-means clustering revealed a total of three subtypes in THCA @ref(fig:THCAhm). This is consistent with the three clusters of THCA observed in the full pan-cancer GSVA data.



**Figure 3.7:** Volcano plot of THCA pathway activity. Downregulated pathways are colored red, upregulated pathways blue. Not significantly altered pathways are colored green. Most significantly altered pathways are labelled with their name.

## RESULTS

The follicular histological type was enriched in cluster B, with no tall cell types present in this cluster. Judging from histological type alone no difference in clusters A and C was observed. Most significant changes in pathway activity were observed in pathways concerning proliferative signaling. In comparison with all other tumor types, cluster A displayed high activity of RAS, JAK/STAT and EWSR1/FL1-fusion mediated signaling as well as elevated signatures associated with carcinogenesis driven by alpha6beta4 activity. In contrast, these pathways were downregulated in cluster B, with it showing elevated activity in mTOR, MAPK, PI3K, and EGFR signaling cascades. Cluster C was found to upregulate all the aforementioned forms of proliferative signaling. All clusters showed a homogenous upregulation of hedgehog, ERBB2, and MST1 pathway activity. Regarding immune response, cluster C showed no significant alterations in the respective hallmark pathways, however, these pathways were downregulated in both clusters A and B. With this data, we can identify two seemingly different forms of proliferative signaling driving carcinogenesis in THCA. These forms can either occur separately as in the case of clusters A and B or combined as for cluster C.



**Figure 3.8:** Pathway activity of the 50 most variant, hallmark, and 20 most significantly altered pathways for each patient. Column clusters were obtained by k-means clustering with k=3. Pathway activities were computed via GSEA of pan-cancer expression data. For all pathway activities see figure (XXX in the appendix).

**THCA subtypes do not differ in their metabolism.**



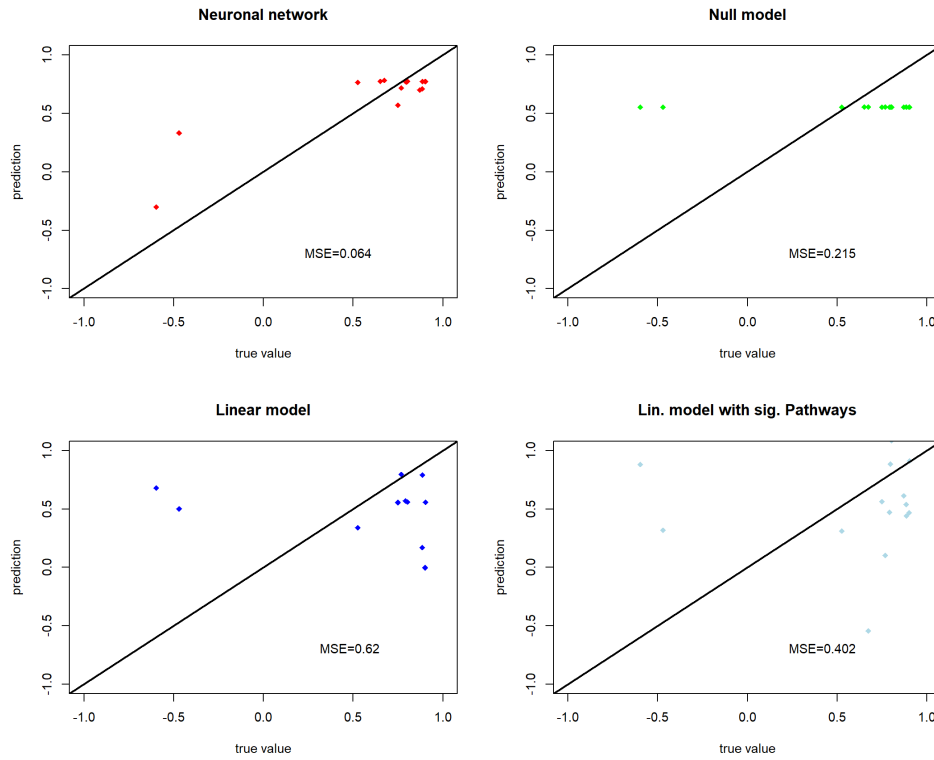
### 3.5 Regression analysis of THCA pathway activity

To select a suitable pathway for regression analysis, the top 20% pathways regarding their variance in activity were chosen, as for the regression model to predict. Pathways with little variance were found to be better predicted by a null model (Fig xxx supplementary material). To factor in biological significance, the intersect of the 25 most significantly altered pathways from GSVA with the high variance pathways was computed. This resulted in three significantly altered and highly variant pathways among which the REACTOME\_INTERLEUKIN\_36\_PATHWAY gene set was selected. This gene set ranks 8th among the highest upregulated pathways with an associated p-value of  $8.411155e-15$ . Interleukin 36 signaling is connected to both MAPK activity and through the activation of NF-kB and also the expression of integrin  $\alpha6\beta4$ . (Bhatia *et al.*, 2013; Liberzon *et al.*, 2015b; Queen *et al.*, 2019).

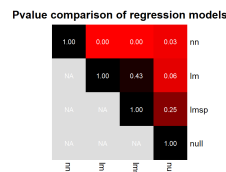
Regression of the REACTOME\_INTERLEUKIN\_36\_PATHWAY gene set showed mixed results. After multiple testing an architecture with two hidden layers with 10 and 20 neurons respectively at ‘`set.seed(50)`’ was shown to produce the best results for neuronal network regression. Among the tested models, the neuronal network performed best on the test data with a mean squared error (MSE) of 0.06. However, the linear regression model failed to predict the data accurately (MSE = 0.62). Repeated linear regression with just pathways contributing significantly to the result the performance was enhanced (MSE = 0.40), however, remained worse than a null model (MSE = +0.22). @ref(fig:reg) Comparison of the MSE on test and training data reveals, that linear model is highly overfitted ( $\Delta\text{MSE} = +0.55$ ) with the linear model with significant pathways fitting slightly better ( $\Delta\text{MSE} = +0.23$ ) to the data. Our null model displayed a good, yet slightly underfitted performance with  $\Delta\text{MSE} = -0.08$ . With a  $\Delta\text{MSE} = +0.009$  the neuronal network shows a perfect fit.

A comparison of the four regression models via the F-test function `var.test()` showed a significant improvement of the neuronal network compared to all other models. All other models showed no significant differences in their performance @ref(fig:reg) compared to each other. From this data, we can conclude that a neuronal network is the best choice for most accurately predicting IL-36 pathway activity in our test data.

## RESULTS



**Figure 3.10:** Regression results for various models on THCA GSEA test data. True values are plotted against predicted values, black slope indicate a perfect prediction.



**Figure 3.11:** F-test comparison of various regression models. p-values are obtained from a two-sided variance test and displayed as heatmap. nn = neuronal network, lm = linear regression, lmsp = linear regression with only significant pathways, null = null model.



## 4 Discussion

Our pan-cancer analysis showed four clusters in pathway activity data. We were able to find specific pathways, which were enriched only in certain histological types like glioblastomas and adenocarcinomas. The focused analysis of THCA expression data revealed pathways driving carcinogenesis in Thyroid cancer. Furthermore we were able to subcategorize THCA into three subtypes based on proliferative signalling pathway activity. As shown by the data signaling through  $\alpha6\beta4$ , RAS, JAK/STAT, and EWSR1/FLI1-fusion mediated pathways are linked to the non-follicular histological subtype of THCA. Our findings from pan-cancer expression data show promising results. Via GSVA analysis we identified four clusters in the cancer types correlating strongly to the associated histological type. Glioblastoma seem to take a special role as they are predominantly characterized by the high activity of neural crest differentiation pathways and receptor tyrosine kinases. This is in line with previous studies showing that glioblastomas derive from neural crest cells (Bednarczyk and McIntyre, 1992).

This was also found for some melanoma like UVM, which explains the observed clustering of UVM with other glioblastoma. Also, the high receptor tyrosine kinase activity has been linked to the formation of UVM and glioblastoma and suggested as a possible target for therapy (Wade *et al.*, 2013; Jo *et al.*, 2019).

Further, especially liver and kidney adenocarcinoma seemed to form a strong subcluster within the other adenocarcinoma. They are characterized by exceptionally high activity of metabolic pathways such as carbohydrate metabolism, lipid, and amino acid synthesis. Again, this change in metabolism was previously found in hepatocellular carcinoma (Sanginetto *et al.*, 2020). An up-regulation of these metabolic pathways may lead to cell growth and proliferation, due to higher metabolic activity, providing more biomass and energy.

The most significant classification we found was the clustering of tumor types by their differentiation stage. Poorly differentiated tumors like leukemia and squamous cell carcinoma show an upregulation of pathways associated with embryonic stem cell-like expression signatures. In contrast highly differentiated tumors like most adenocarcinoma as well as most glioblastoma underexpress these gene sets. Such a clustering by differentiation stage was previously described

by Ben-Porath *et al.*. However, these findings cannot be verified directly as provided annotation data did not contain information regarding the differentiation stage (Ben-Porath *et al.*, 2008).

Taken together our results are in line with current research and allow for the following hypothesis: The expression profile of a given cancer type depends highly on its differentiation stage and its histological type but little on the actual tumor type itself. Understanding how these changes in expression link to mutational signatures might help in developing druggable targets for therapy.

From our GSEA and pan-cancer GSVA results, we identify two separate ways of carcinogenesis in THCA. The follicular subtype upregulates proliferative signaling through mTOR/PI3K and MAPK signaling pathways. This was previously shown by Furuya *et al.* (Furuya *et al.*, 2007).

A second way of carcinogenesis by signaling through alpha6beta4, RAS, JAK/STAT, and EWSR1/FLI1-fusion mediated pathways was observed in the data. This way of carcinogenesis was linked to non-follicular types of THCA (Noh *et al.*, 2010; Oliveira *et al.*, 2017; Bi *et al.*, 2019).

The finding that the alpha6beta4 integrin signaling pathway and associated pathways such as IL-36 signaling and Typ I hemidesmosome synthesis were significantly enhanced in THCA is in line with previous studies. Those studies link alpha6beta4 signaling to the development of aggressive forms of thyroid cancer Noh *et al.* (2010). Also, oncogenic signaling pathways commonly associated with different cancer types were significantly upregulated in THCAs. Among them, we observed ERBB2 and MST1 signaling commonly found in breast cancer. A role for MSP/Ron in breast cancer has recently been elucidated, wherein this pathway regulates tumor growth, angiogenesis and metastasis (Kretschmann *et al.*, 2010).

A main pathway up-regulated in THCA is signaling through EWSR1/FLI1-fusion protein, while non-histone protein methylation is down-regulated in THCA. This process was identified as an import modulator of intracellular signaling by the MAPK, WNT, BMP, Hippo, and JAK/STAT pathways and might play an important role as a driver of carcinogenesis in THCA (Biggar and Li, 2015). Together these findings give a general overview of mechanisms driving carcinogenesis in THCA. However, no information about possible THCA subtypes or differences in pathway activity between patients can be obtained from this data.

Pan-cancer GSVA shows three distinct clusters in the expression data, upregulating either one or both ways of proliferative signaling. While the follicular subtype seemed to strongly correlate with one cluster, a similar process was not observed in tall-cell and classical phenotypes. With more detailed annotation data it might be possible to link anaplastic and papillary histological subtypes of THCA to the two yet unassigned clusters.

Despite differences in proliferative signaling, all clusters share an upregulated hedgehog signaling pathway which is consistent with the literature (Hinterseher *et al.*, 2014). Also, metabolic changes in line with the Warburg effect were observed in all clusters.

For regression analysis by linear regression and a neural network the REACTOME\_INTERLEUKIN\_36\_PATHWAY was chosen, since it is connected to both MAPK activity and through the activation of NF-kB and also the expression of integrin alpha6beta4. An effective regression might be crucial in finding potentially druggable targets in combating THCA. Our data suggest that our neuronal network is well suited to predict pathway activities from GSEA data. The model shows an excellent fit to data and produces only minor errors. However, both linear models struggle in predicting the data accurately. This might be since GSEA pathway activity data usually clusters into an up- and down-regulated group with no values in between. Since the REACTOME\_INTERLEUKIN\_36\_PATHWAY also shows this problem, the two clusters might produce larger correlation values that might impact the accuracy of the regression coefficients and the intercept. Secondly, the correlation of the residuals with the test data values did not approach zero, thus, our linearity assumption is not met. Therefore, it can be concluded that a linear regression model is not well suited to predict the REACTOME\_INTERLEUKIN\_36\_PATHWAY activity accurately.

## 4.1 Conclusion

xxx fehlt noch und kommt dann zur diskussion

## 5 Outlook

Futher ways of analysis could be the prediction of the histological type of THCA as well as the way of carcinogenesis with a neuronal network. This might be possible with a larger training data set as well as more detailed and specified annotations. Furthermore, it might be possible to link whole genome sequencing and methylation data to pathway activity. In that way, one could suggest a suitable targeted therapy option for a THCA patient based only on sequencing data from a small biopsy sample.

## 6 References

- Alberts, J, B., and Walter, P (2015). Molecular biology of the cell, New York: Garland science.
- Bednarczyk, JL, and McIntyre, BW (1992). Expression and ligand-binding function of the integrin  $\alpha 4 \beta 1$  (VLA-4) on neural-crest-derived tumor cell lines. *Clinical & Experimental Metastasis* 10, 281–290.
- Ben-Porath, I, Thomson, MW, Carey, VJ, Ge, R, Bell, GW, Regev, A, and Weinberg, RA (2008). An embryonic stem cell-like gene expression signature in poorly differentiated aggressive human tumors. *Nat Genet* 40, 499–507.
- Bhatia, V, Mula, RV, and Falzon, M (2013). Parathyroid hormone-related protein regulates integrin  $\alpha 6$  and  $\beta 4$  levels via transcriptional and post-translational pathways. *Exp Cell Res* 319, 1419–1430.
- Bi, C-L, Zhang, Y-Q, Li, B, Guo, M, and Fu, Y-L (2019). Retracted: MicroRNA-520a-3p suppresses epithelial–mesenchymal transition, invasion, and migration of papillary thyroid carcinoma cells via the JAK1-mediated JAK/STAT signaling pathway. *Journal of Cellular Physiology* 234, 4054–4067.
- Biggar, KK, and Li, SSC (2015). Non-histone protein methylation as a regulator of cellular signalling and function. *Nature Reviews Molecular Cell Biology* 16, 5–17.
- Cabanillas, ME, McFadden, DG, and Durante, C (2016). Thyroid cancer. *Lancet* 388, 2783–2795.
- Coca-Pelaz, A et al. (2020). Papillary thyroid cancer-aggressive variants and impact on management: A narrative review. *Adv Ther* 37, 3112–3128.
- Durinck, S, Spellman, PT, Birney, E, and Huber, W (2009). Mapping identifiers for the integration of genomic datasets with the r/bioconductor package biomaRt. *Nature Protocols* 4, 1184–1191.
- Fritsch, S, Guenther, F, and Wright, MN (2019). Neuralnet: Training of neural networks.
- Furuya, F, Lu, C, Willingham, MC, and Cheng, S (2007). Inhibition of phosphatidylinositol 3-kinase delays tumor progression and blocks metastatic spread in a mouse model of thyroid cancer. *Carcinogenesis* 28, 2451–2458.
- Hanahan, D, and Weinberg, RA (2011). Hallmarks of cancer: The next generation. *Cell* 144, 646–674.

- Hänzelmann, S, Castelo, R, and Guinney, J (2013a). GSVA: Gene set variation analysis for microarray and RNA-seq data. *BMC Bioinformatics* 14, 7.
- Hänzelmann, S, Castelo, R, and Guinney, J (2013b). GSVA: Gene set variation analysis for microarray and RNA-Seq data. *BMC Bioinformatics*.
- Hao, Y et al. (2021). Integrated analysis of multimodal single-cell data. *Cell* 184, 3573–3587.e29.
- Hinterseher, U, Wunderlich, A, Roth, S, Ramaswamy, A, Bartsch, DK, Hauptmann, S, Greene, BH, Fendrich, V, and Hoffmann, S (2014). Expression of hedgehog signalling pathway in anaplastic thyroid cancer. *Endocrine* 45, 439–447.
- Jo, DH, Kim, JH, and Kim, JH (2019). Targeting tyrosine kinases for treatment of ocular tumors. *Archives of Pharmacal Research* 42, 305–318.
- Kant, R, Davis, A, and Verma, V (2020). Thyroid nodules: Advances in evaluation and management. *Am Fam Physician* 102, 298–304.
- Konopka, T (2022). Umap: Uniform manifold approximation and projection.
- Korotkevich, G, Sukhov, V, and Sergushichev, A (2019). Fast gene set enrichment analysis. *bioRxiv*.
- Kretschmann, KL, Eyob, H, Buys, SS, and Welm, AL (2010). The macrophage stimulating protein/ron pathway as a potential therapeutic target to impede multiple mechanisms involved in breast cancer progression. *Curr Drug Targets* 11, 1157–1168.
- Liberzon, A, Birger, C, Thorvaldsdóttir, H, Ghandi, M, Mesirov, JP, and Tamayo, P (2015b). The molecular signatures database (MSigDB) hallmark gene set collection. *Cell Syst* 1, 417–425.
- Liberzon, A, Birger, C, Thorvaldsdóttir, H, Ghandi, M, Mesirov, JP, and Tamayo, P (2015a). The molecular signatures database (MSigDB) hallmark gene set collection. *Cell Syst* 1, 417–425.
- Lin, JD (2007). Papillary thyroid carcinoma with lymph node metastases. *Growth Factors* 25, 41–49.
- Lunt, M (2013). Introduction to statistical modelling: Linear regression. *Rheumatology* 54, 1137–1140.
- Maaten, L van der, and Hinton, G (2008). Visualizing data using t-SNE. *Journal of Machine Learning Research* 9, 2579–2605.
- Noh, TW, Soung, YH, Kim, HI, Gil, HJ, Kim, JM, Lee, EJ, and Chung, J (2010). Effect of beta4 integrin knockdown by RNA interference in anaplastic thyroid carcinoma. *Anticancer Res* 30, 4485–4492.
- Oliveira, G, Polónia, A, Cameselle-Teijeiro, JM, Leitão, D, Sapia, S, Sobrinho-Simões, M, and Eloy, C (2017). EWSR1 rearrangement is a frequent event in papillary thyroid carcinoma and

## REFERENCES

---

- in carcinoma of the thyroid with ewing family tumor elements (CEFTE). *Virchows Archiv* 470, 517–525.
- Prete, A, Borges de Souza, P, Censi, S, Muzza, M, Nucci, N, and Sponziello, M (2020). Update on fundamental mechanisms of thyroid cancer. *Front Endocrinol (Lausanne)* 11, 102.
- Queen, D, Ediriweera, C, and Liu, L (2019). Function and regulation of IL-36 signaling in inflammatory diseases and cancer development. *Front Cell Dev Biol* 7, 317.
- Rabinovitz, I, and Mercurio, AM (1996). The integrin alpha 6 beta 4 and the biology of carcinoma. *Biochem Cell Biol* 74, 811–821.
- Reimand, J et al. (2019). Pathway enrichment analysis and visualization of omics data using g:profiler, GSEA, cytoscape and EnrichmentMap. *Nature Protocols* 14, 482–517.
- Riedmiller, MA Rprop - description and implementation details.
- Sanginetto, M, Villani, R, Cavallone, F, Romano, A, Loizzi, D, and Serviddio, G (2020). Lipid metabolism in development and progression of hepatocellular carcinoma. *Cancers* 12, 1419.
- Sharma, S, Quinn, D, Melenhorst, JJ, and Pruteanu-Malinici, I (2021). High-dimensional immune monitoring for chimeric antigen receptor t cell therapies. *Current Hematologic Malignancy Reports* 16, 112–116.
- Tsibulnikov, S, Maslov, L, Voronkov, N, and Oeltgen, P (2020). Thyroid hormones and the mechanisms of adaptation to cold. *Hormones (Athens)* 19, 329–339.
- Wade, A, Robinson, AE, Engler, JR, Petritsch, C, James, CD, and Phillips, JJ (2013). Proteoglycans and their roles in brain cancer. *The FEBS Journal* 280, 2399–2417.
- Wang, X, Pei, Z, Hossain, A, Bai, Y, and Chen, G (2021). Transcription factor-based gene therapy to treat glioblastoma through direct neuronal conversion. *Cancer Biol Med* 18, 860–874.

## 7 Appendix