



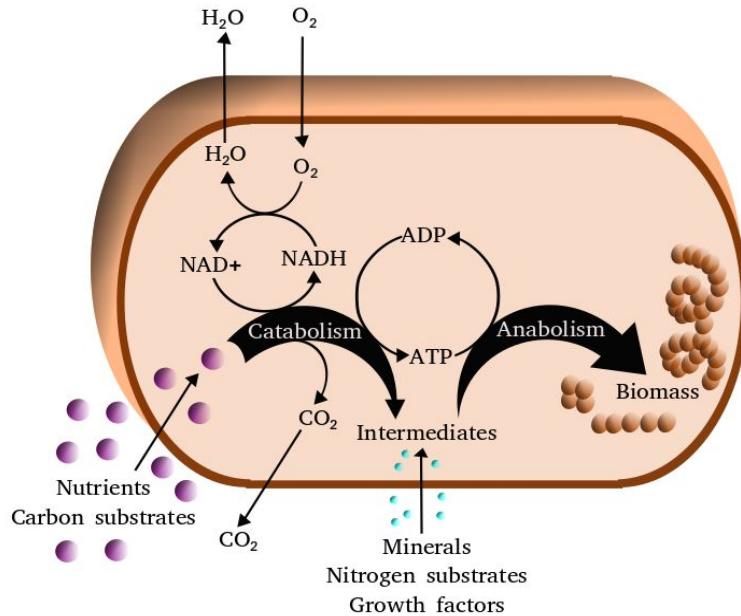
Gene expression is a poor predictor of the metabolite abundance in cancer cells

Huaping Li
School of Biomedical Sciences, HKU
2022.03

Background

Metabolism refers to all biochemical processes taking place in live cells and organisms.

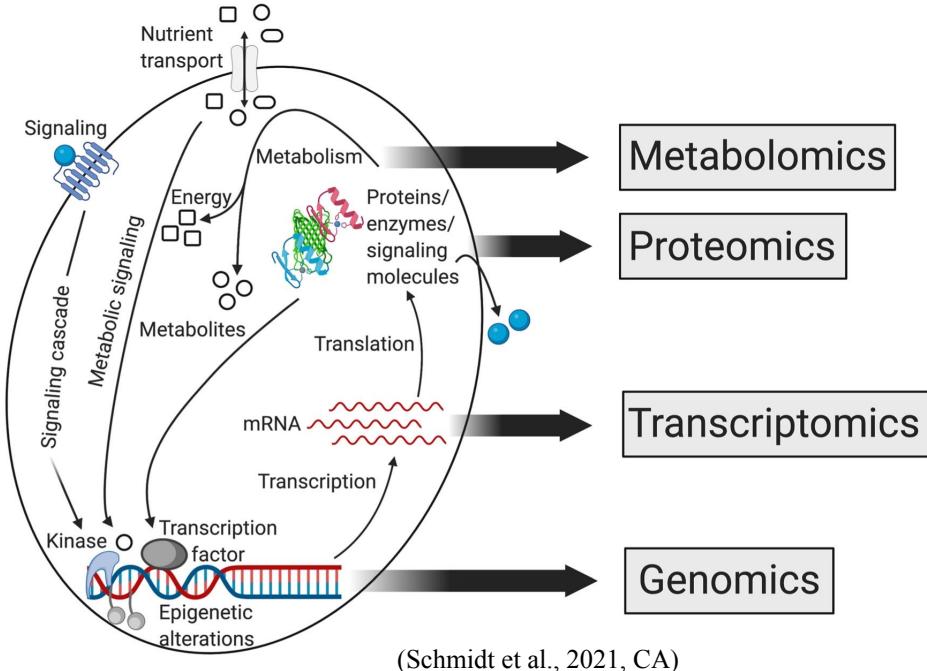
Metabolites are the by-products or products of cellular metabolic activities and are usually shown as small molecules, such as amino acids, lactic acid, lipids, nucleotide, vitamins, ethanol, and glycerol, etc.



Simplified view of the cellular metabolism

Background

- Cancer cells tend to modulate their metabolism to meet the high anabolic requirements of rapid proliferation, and therefore promotes cancer progression(DeBerardinis et al., 2007; Liberti & Locasale, 2016)
- The global changes in metabolite levels has been reported as a hallmark of cancer(Faubert et al., 2020).



Research Aim

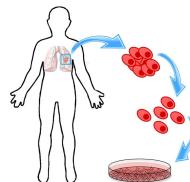
Metabolomics
less resources
hard to capture & measure

Transcriptomics
huge resources
easily detect & profile

Database

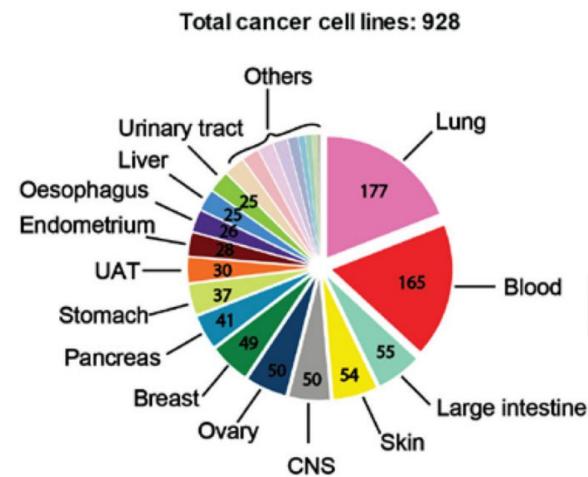
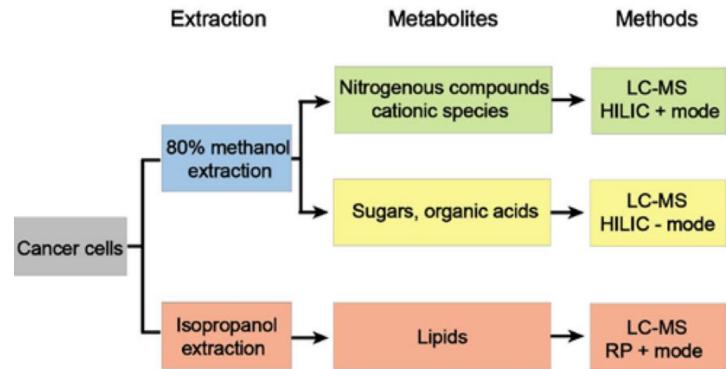


large-scale genetic characterization of
~1000 cancer cell lines



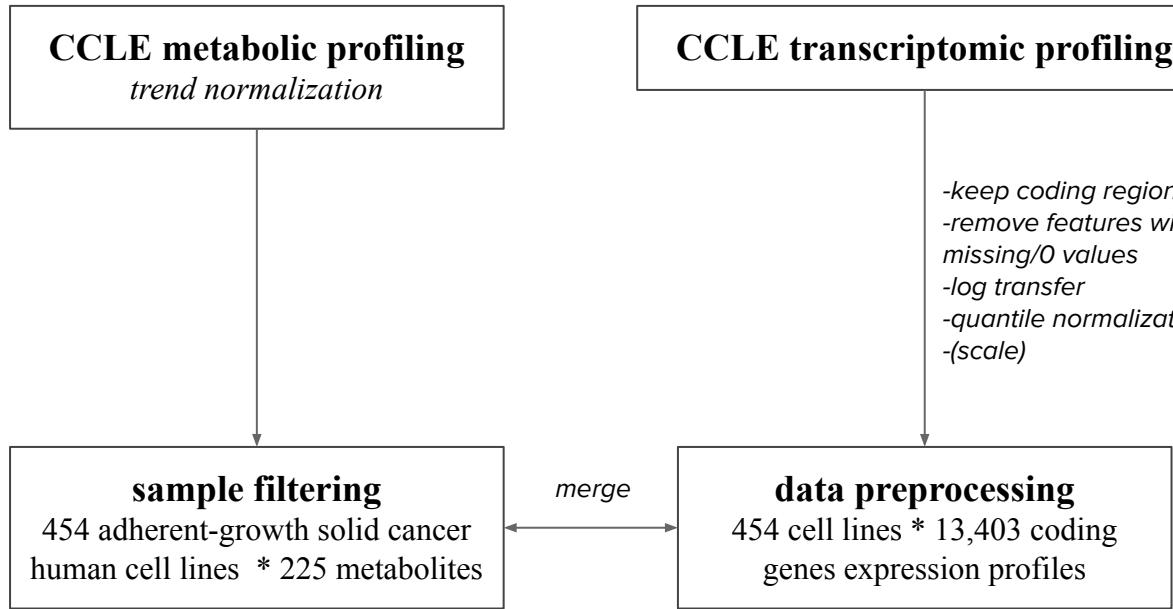
Metabolomics:
225 metabolites * 928 cell lines (29 cancer types)

Transcriptomics:
28k genes * 1019 cell lines

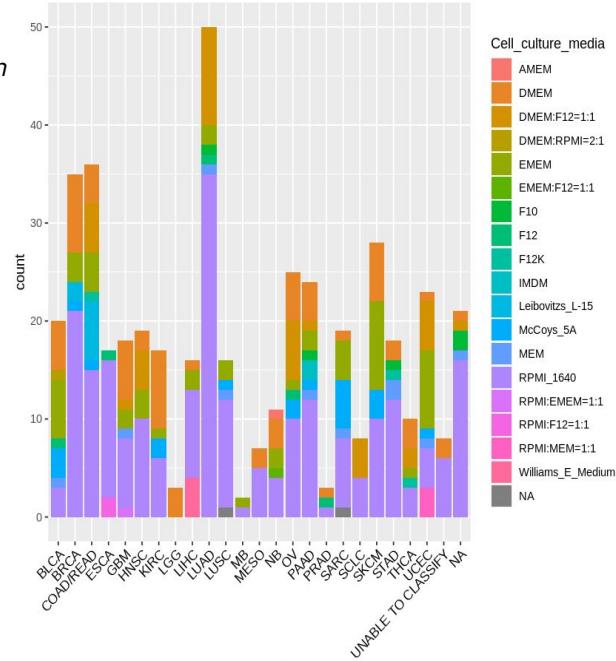


(Li et al., 2019, Nature Medicine)

Analysis pipeline

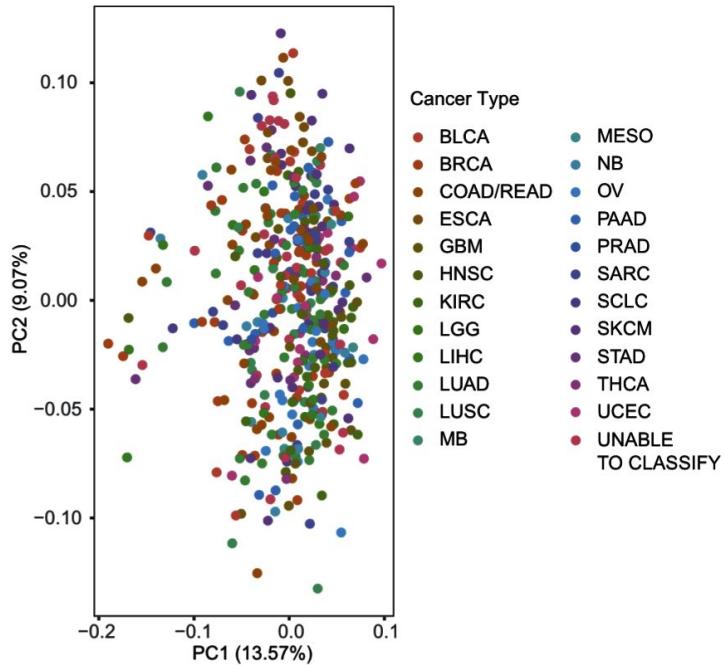
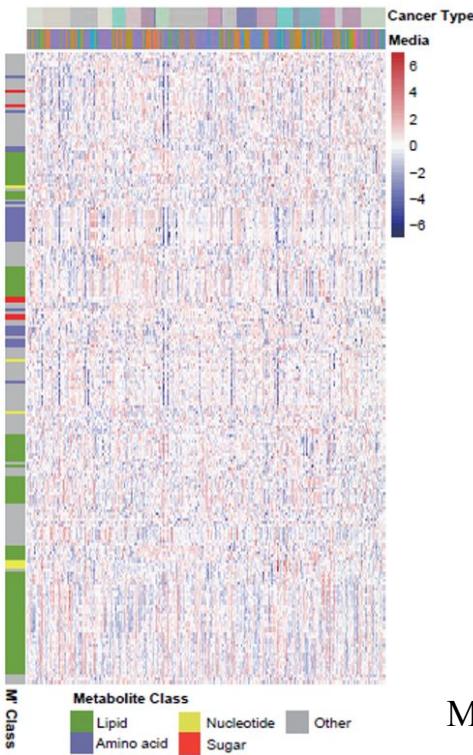


Basic statistics analysis:

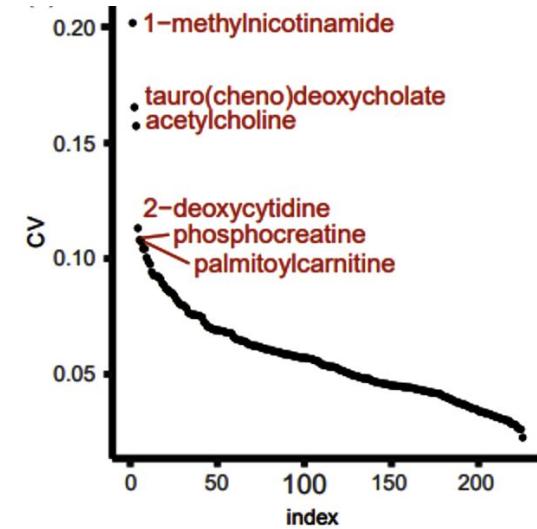


Metabolism heterogeneity

- “Metabolic phenotypes in tumors are both heterogeneous and flexible” (Kim & DeBerardinis, 2019)

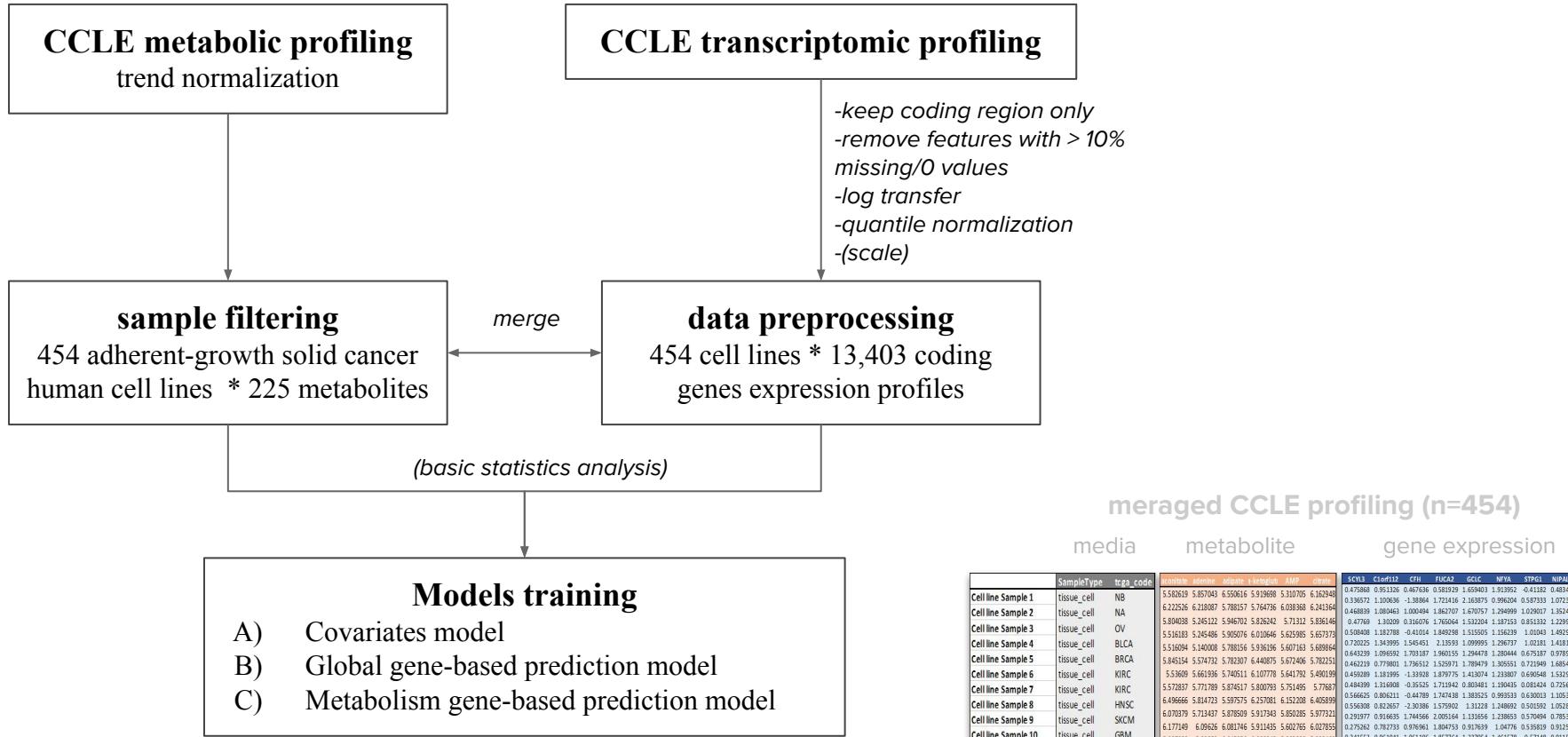


Metabolism heterogeneity among cancer cell lines



Coefficient of variation(CV) of each metabolite.[$CV = SD/\text{mean}$]

Analysis pipeline



Linear & Lasso regression

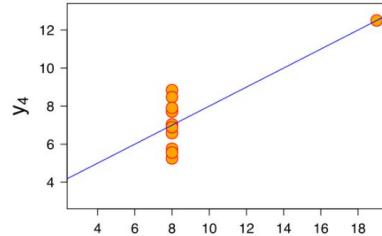
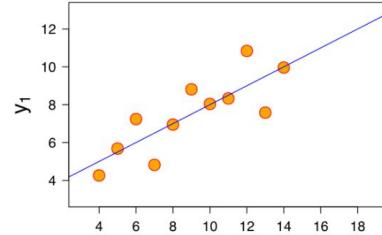
“In statistics, **linear regression** is a linear approach for **modelling the relationship** between a scalar response and one or more explanatory variables”

$$y \sim w_1x_1 + w_2x_2 + \dots + w_nx_n$$

multicollinearity (eg. $x_1=\log x_2$)
feature bias

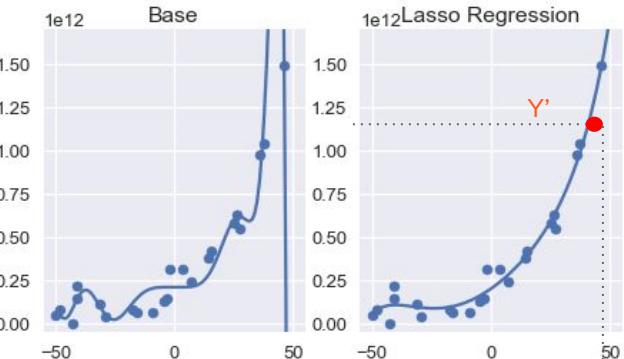
Cellline Sample	y
Cellline Sample 1	3.553519
Cellline Sample 2	6.222526
Cellline Sample 3	5.884038
Cellline Sample 4	5.532257
Cellline Sample 5	5.532624
Cellline Sample 6	5.881514
Cellline Sample 7	5.533693
Cellline Sample 8	6.006666
Cellline Sample 9	6.073379
Cellline Sample 10	6.177149

SCGB3	Claud13	CFL	FUCA2	GCF	MTH	SPD5	NPM013	LASS1	ENPPA
0.475688	0.051326	0.467936	0.581529	1.659463	0.313152	0.413132	0.485405	1.847922	0.216263
0.433432	0.051326	0.467936	0.581529	1.659463	0.313152	0.413132	0.485405	1.847922	0.216263
0.468819	1.000441	1.000494	1.362707	1.670577	1.249499	1.029327	1.362484	1.516034	0.518077
0.47769	1.362029	1.362066	1.768654	1.532104	1.375153	0.815132	1.2099	1.532072	0.5867
0.500000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
0.720225	1.343995	1.464571	1.235030	1.099955	1.287337	1.021381	1.418132	1.790572	0.5867
0.643139	1.096392	1.703137	1.460155	1.284678	1.265644	0.673137	0.739108	1.775208	0.846202
0.464262	0.051326	0.467936	0.581529	1.659463	0.313152	0.413132	0.485405	1.847922	0.216263
0.439384	1.181071	1.370975	1.430701	1.238007	1.099948	1.520204	1.410933	0.770702	0.000000
0.448443	0.051326	0.467936	0.581529	1.659463	0.313152	0.413132	0.485405	1.847922	0.216263
0.564625	0.802311	0.647789	1.747631	1.035252	0.993533	0.430033	1.105458	1.8994	0.43103
0.586808	0.922657	2.031707	1.579002	1.31228	1.248902	0.502032	0.826507	1.850909	0.786244
0.523930	0.051326	0.467936	0.581529	1.659463	0.313152	0.413132	0.485405	1.847922	0.216263
0.757582	0.782733	0.179763	1.040753	0.917039	1.047376	0.538018	0.912507	1.506909	0.021074
0.341553	0.961341	1.061196	1.857761	1.377654	1.451178	0.71748	0.911518	1.85624	0.05105



“**Lasso** is a modification of linear regression, where the model is penalized for the sum of absolute values of the weights. It performs both **variable selection** and **regularization** in order to enhance the prediction accuracy and interpretability of the resulting statistical model.”

$$\frac{1}{2m} \sum_{i=1}^m (y - Xw)^2 + \text{alpha} \sum_{j=1}^p |w_j| \quad \text{penalize}$$

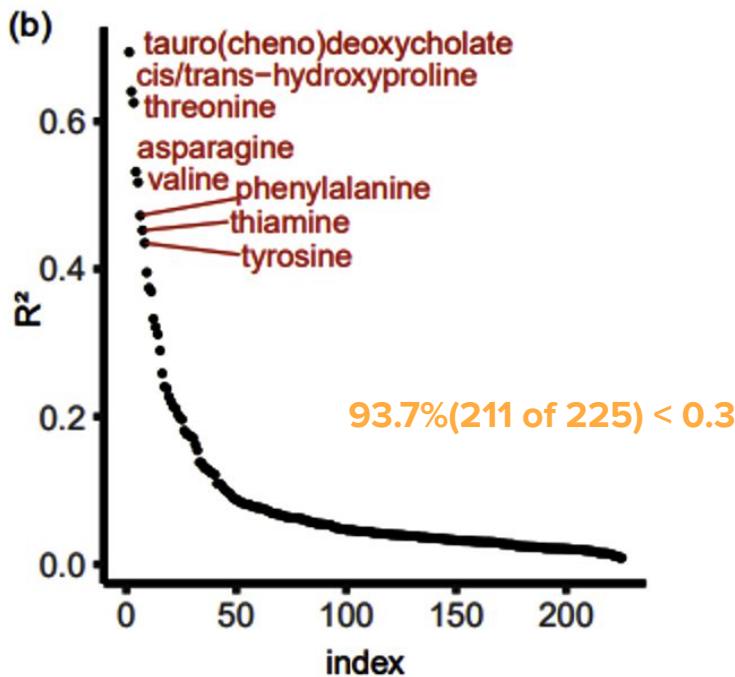


Covariates analysis

10-fold cv, step-wise, linear regression models to check the influence from media and clinical traits

lm: $y \sim 55$ medium componunts + clinical information

(inferred ethnicity, patient age, pathology, cancer type, mutation rates and doubling time (hr))

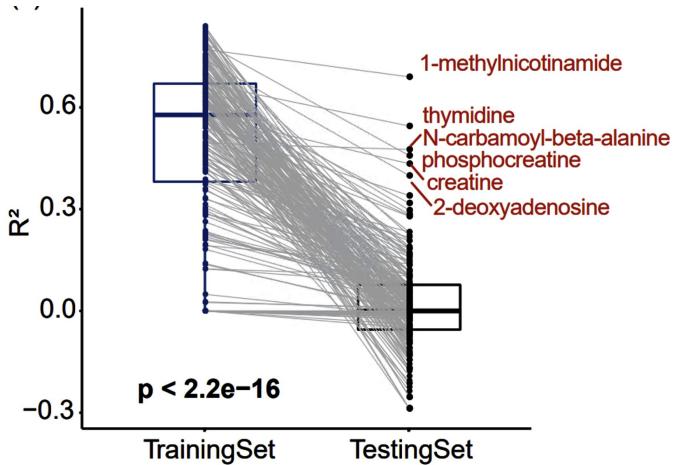
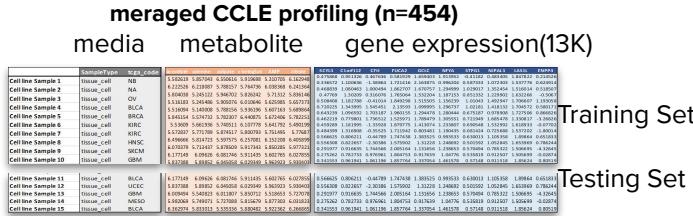


Metabolite	R2
taurodeoxycholate/taurochenodeoxycholate	0.693
cis/trans-hydroxyproline	0.64
threonine	0.625
asparagine	0.531
valine	0.517
phenylalanine	0.472
thiamine	0.452
tyrosine	0.435
lysine	0.395
histidine	0.374
isoleucine	0.369
arginine	0.332
niacinamide	0.321
tryptophan	0.311

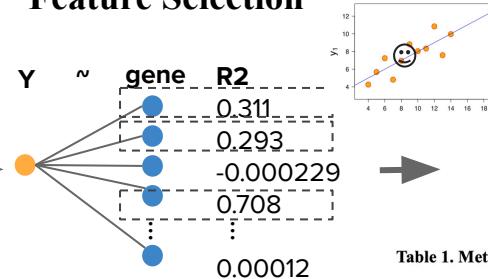
Global gene-based prediction models

linear model: $y \sim \text{cancer type} + \text{media type} + \text{gene}$

Lasso model: $y \sim \text{selected genes}$



Feature Selection



Model training

$Y \sim \text{Top 10\% high } R^2 \text{ genes of each metabolite as features}$

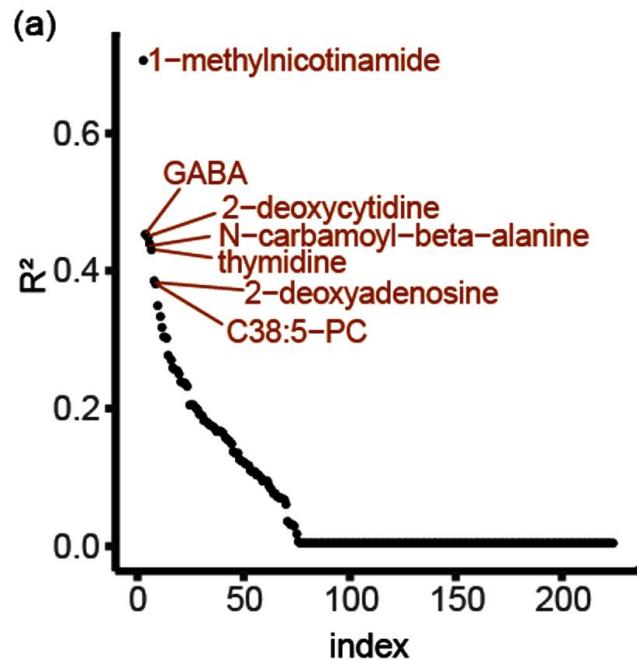
~ ● + ● + ... + ●

Table 1. Metabolites with high R-squared(R^2) from the prediction models.

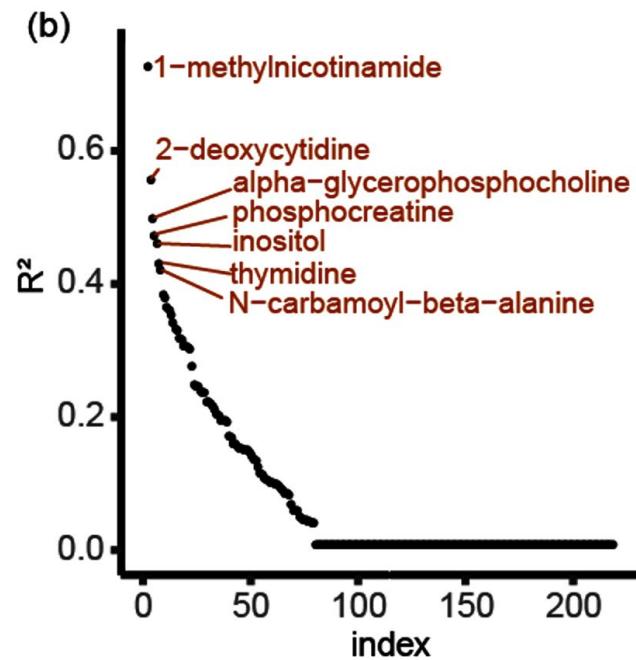
	Global gene-based LASSO model		covariates linear model [▲]
Metabolite	TrainingSet_R ²	TestingSet_R ²	R ²
1-methylnicotinamide	0.771	0.69	0.0612
thymidine	0.67	0.545	0.0622
N-carbamoyl-beta-alanine	0.495	0.476	0.0625
Phosphocreatine	0.79	0.458	0.0217
creatine	0.575	0.434	0.0479
2-deoxyadenosine	0.741	0.399	0.0395

Metabolism gene-based prediction models

Lasso regression: $Y \sim$ all metabolism genes



554 metabolism genes from **SMPDB**



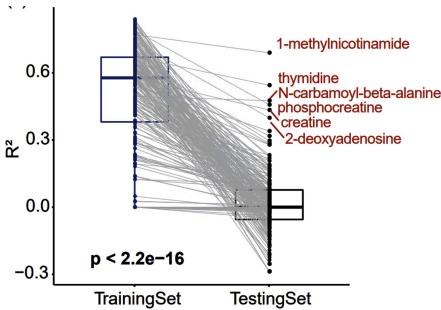
1701 metabolism genes **Reactome**

Summary

Gene expression alone cannot be used to estimate metabolism activity

- Even when accounting for cell culture conditions or cell lineage in the model, few metabolites could be accurately predicted.

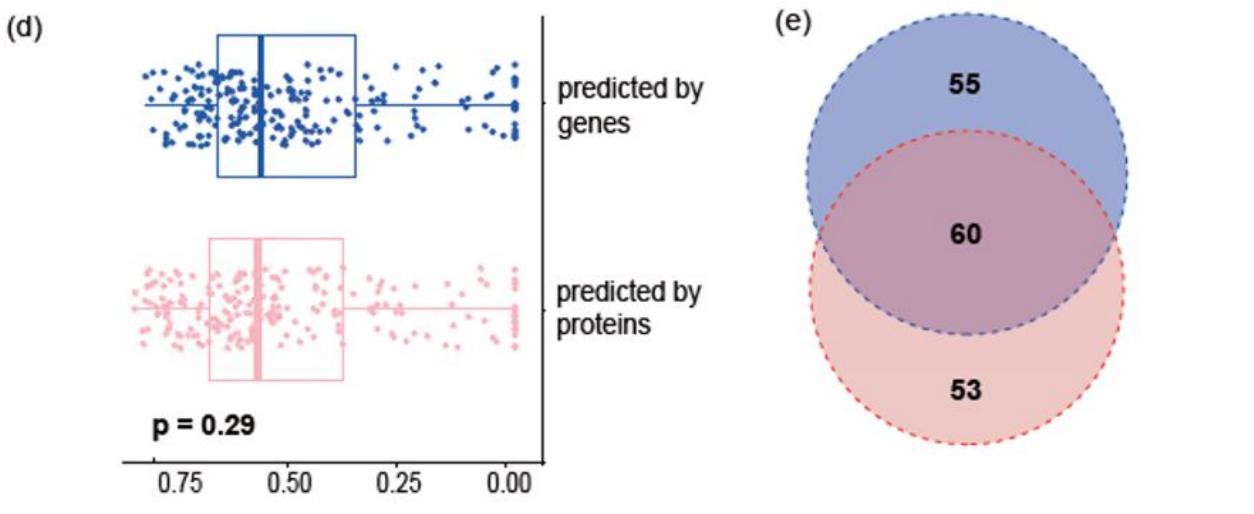
Validated in proteomics level



Global gene-based prediction models



Proteomics-based prediction models



d) R^2 distribution of global gene-based Lasso models (blue) and proteomics-based Lasso models (red). e) Venn diagram of common predictable metabolites (Top 50% high R^2 metabolites in training set) between global gene-based (blue) and proteomics-based (red) Lasso models.

Using associated metabolites in prediction

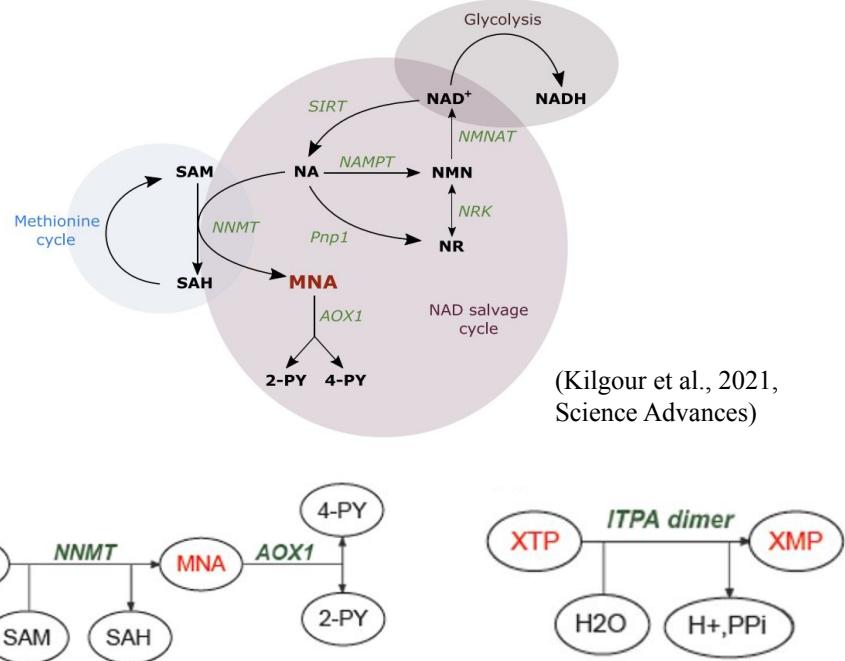


Table 2. Prediction ability change with reactant metabolite level

Prediction Model	R ² in testing set	Reactant Coefficient(rank)
MNA ~ selected genes	0.683	\
MNA ~ selected genes + NAM	0.683	\
XMP ~ selected genes	-0.422	\
XMP ~ selected genes + XTP	0.376	0.57(1)
phosphocreatine ~ selected genes	0.516	\
phosphocreatine ~ selected genes + creatine	0.646	0.902(1)
thymidine ~ selected genes	0.525	\
thymidine ~ selected genes + thymine	0.649	0.399(1)

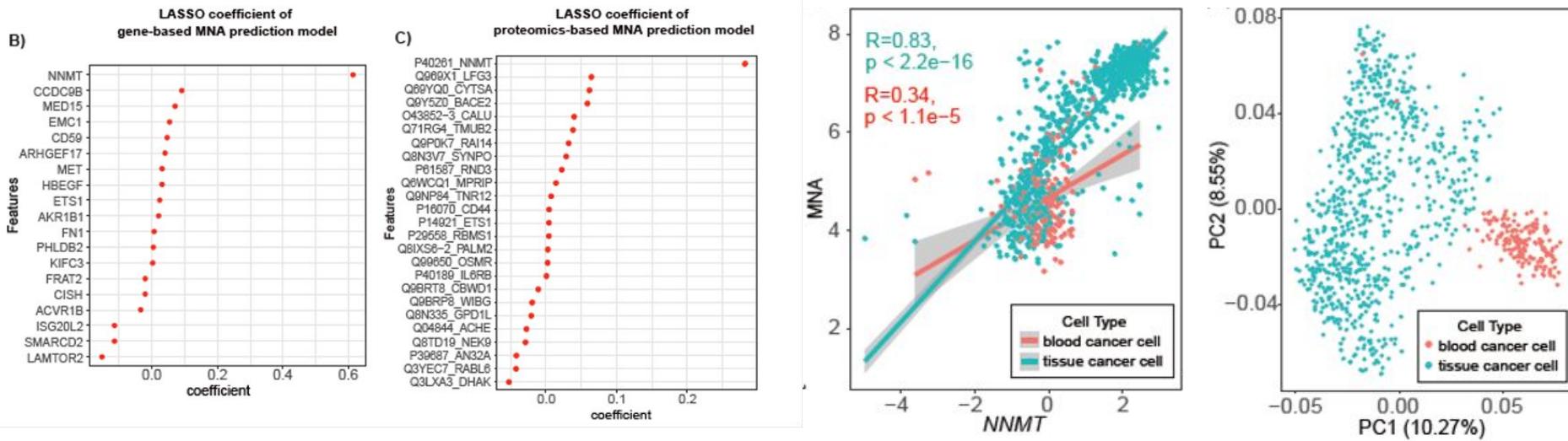
MNA: 1-methylnicotinamide; NAM: Nicotinamide; XMP: xanthosine; XTP: xanthine.

Summary

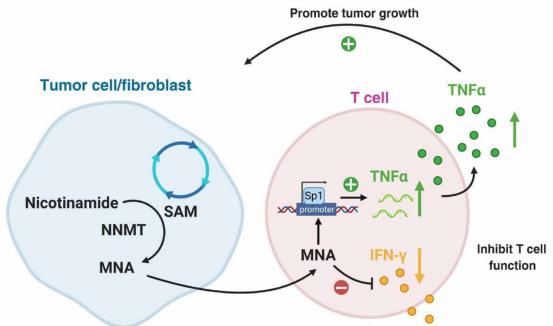
Gene expression alone cannot be used to estimate metabolism activity

- Even when accounting for cell culture conditions or cell lineage in the model, few metabolites could be accurately predicted.
- In some cases, the inclusion of the upstream and downstream metabolites are needed in prediction

NNMT correlated well with MNA in CCLE cohort



1. Solid cancer cells and blood cancer cells show substantially different metabolic profiles
2. Some metabolites maybe cancer-type specific
 - considering cell type and even cancer purity in metabolites prediction.



(Kilgour et al., 2021,
Science Advances)

Summary

Gene expression alone cannot be used to estimate metabolism activity

- Even when accounting for cell culture conditions or cell lineage in the model, few metabolites could be accurately predicted.
- In some cases, the inclusion of the upstream and downstream metabolites are needed in prediction
- The metabolism heterogeneity and sample purity enlarge the challenge of metabolite prediction.

Summary

Gene expression alone cannot be used to estimate metabolism activity

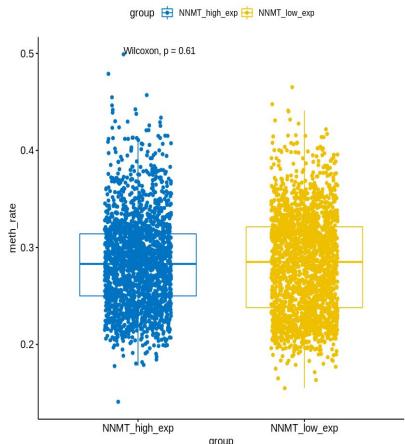
- Even when accounting for cell culture conditions or cell lineage in the model, few metabolites could be accurately predicted.
- In some cases, the inclusion of the upstream and downstream metabolites are needed in prediction
- The metabolism heterogeneity and sample purity enlarge the challenge of metabolite prediction.

Why do most of the prediction models don't perform well?

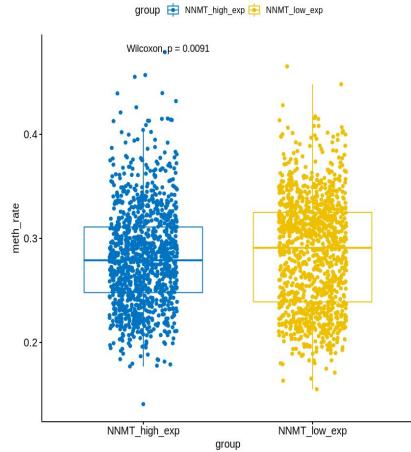
- Gene expression may not reflect the activity of some metabolism associated enzymes. The enzyme activities could be regulated allosterically and post-translationally;
- The bi-directional effect of metabolism and gene expression, making the metabolite level much more unpredictable.
- There are many flexible transcriptional solutions to regulate the metabolism activities.

TCGA bulk samples - NNMT vs methy rate,purity

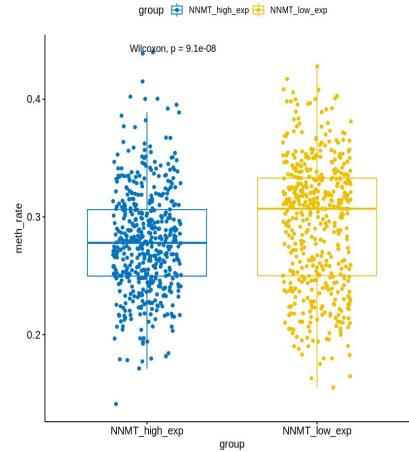
150k detected “Island” methylation sites



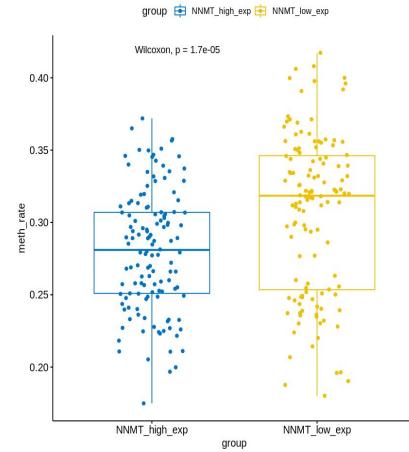
purity = 0.8,n=4k+



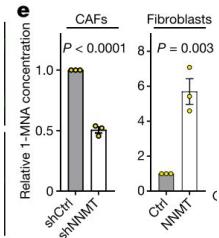
purity = 0.90,n=2k+



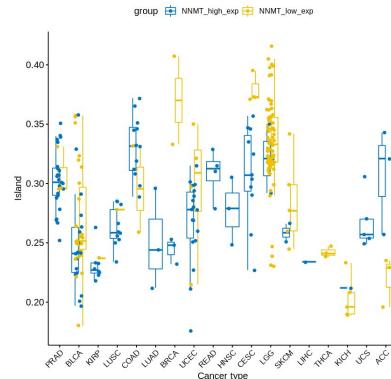
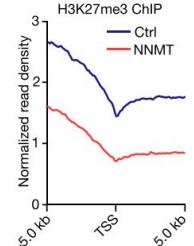
purity = 0.95,n=1k+



purity = 0.98,n=200+



(eckert et al., 2019, Nature)



Future Plans

- 1. Study the possible mechanism of methylation regulation by MNA and NNMT in tumors.**
- 2. See if there are any metabolism signals associated with anti-cancer drug sensitivity.**

Pharmacologic data:

4686 drugs * 579 cell lines (Corsello et al., 2020, Nat Cancer)
264 anti-cancer drugs * 1001 cell lines (Iorio et al., 2016, Cell)

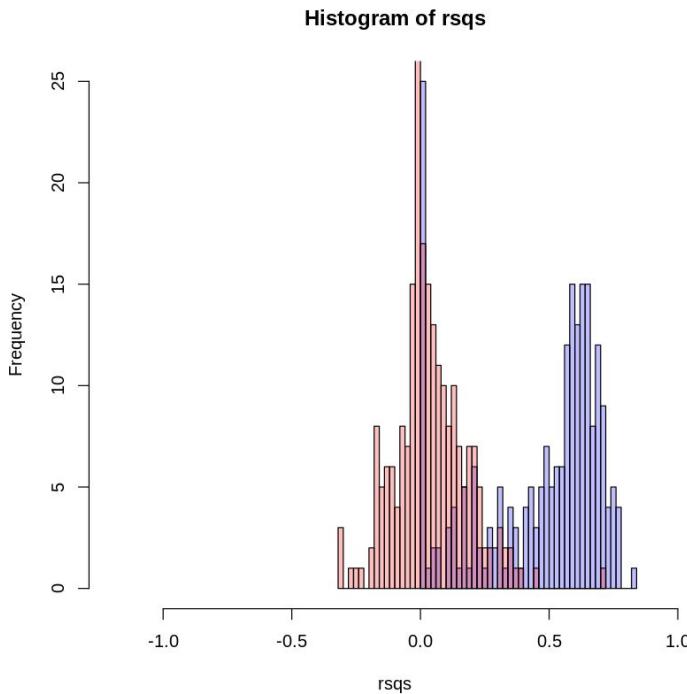
THANKS

bk pages

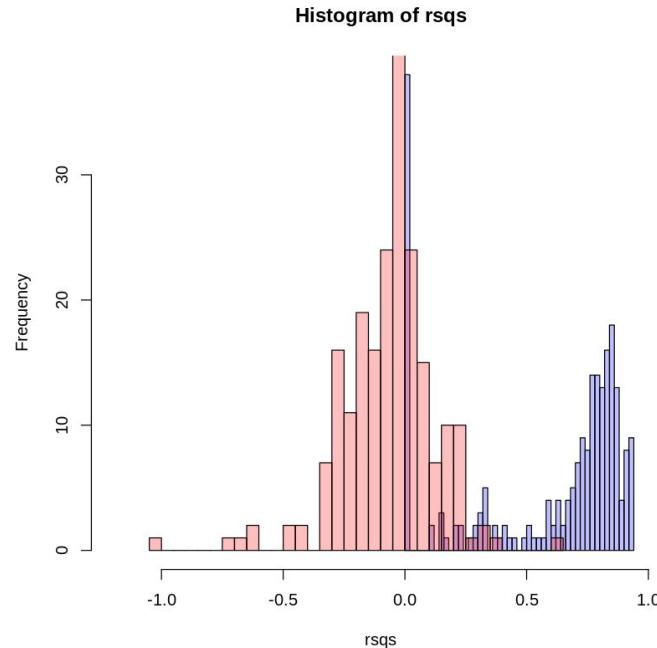
double check

majority of the metaboilites are not influenced by common cell culture media.

multiple media , normed
cancer cell lines (**n=454**)



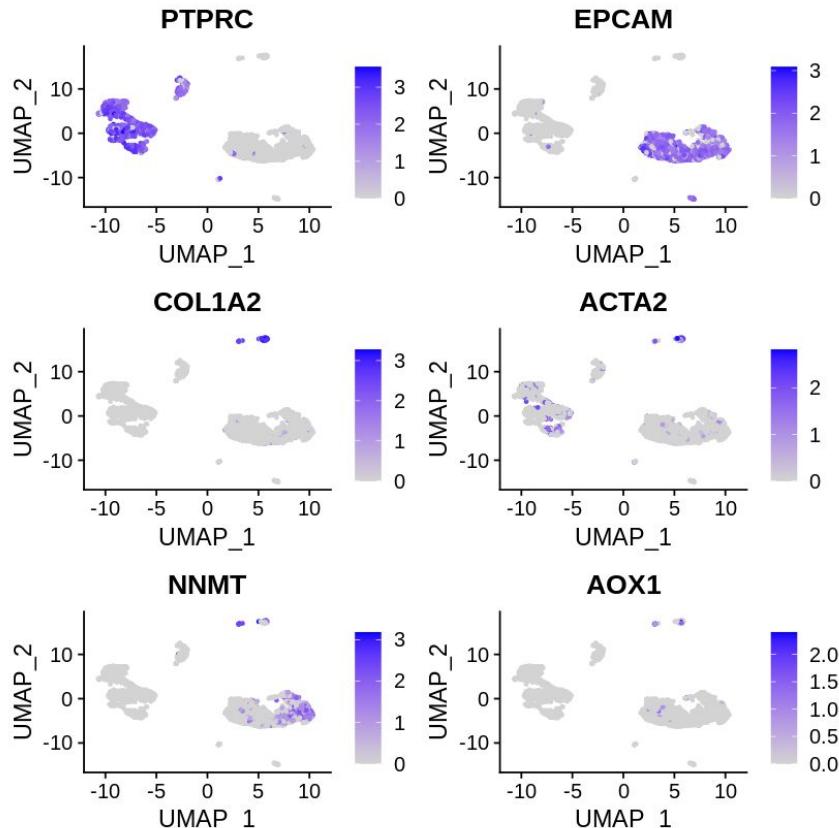
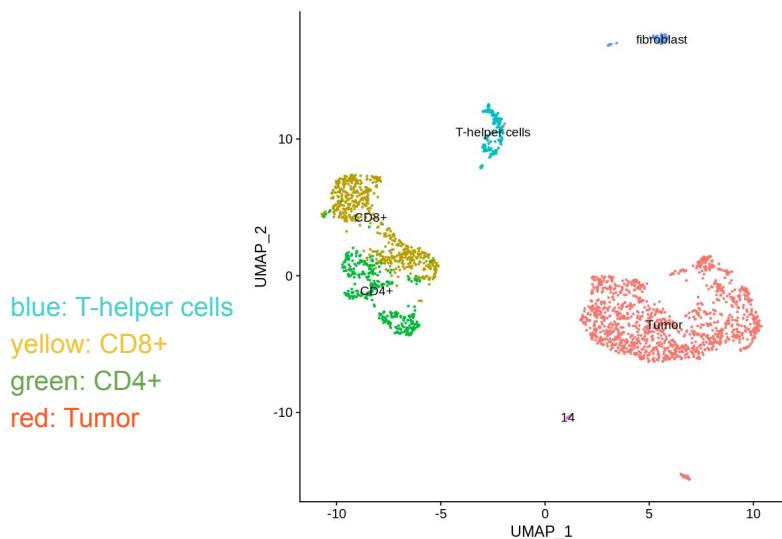
RPMI_1640 ,normed
cancer cell lines (**n=226**)



ovarian Single cell cohort - NNMT distribution in cells

"EpCAM+ (tumor) cells"

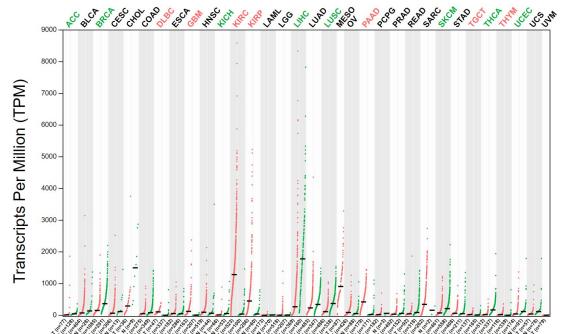
EPCAM (pan-epithelial cell marker),
COL1A2 and ACTA2 (fibroblast markers)
PTPRC (CD45+, immune cell marker)



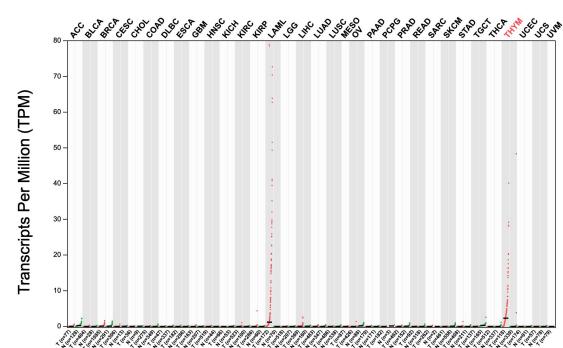
metabolism pathway - gene and metabolite distribution among cancer types

1. NNMT expressed in tumor cell, while most AOX1 not
2. AOX1 highly expressed in LAML tumor cells
3. MNA is not predictable in LAML samples

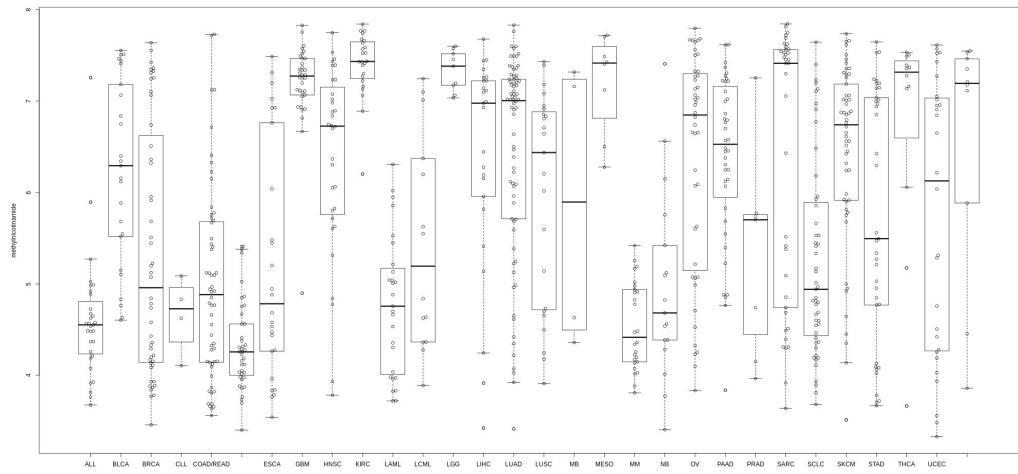
NNMT



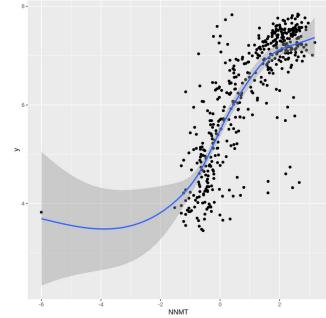
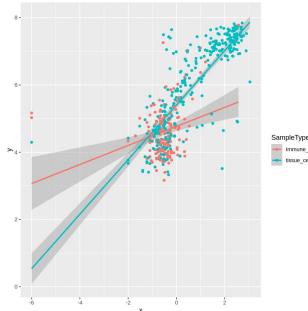
AOX1



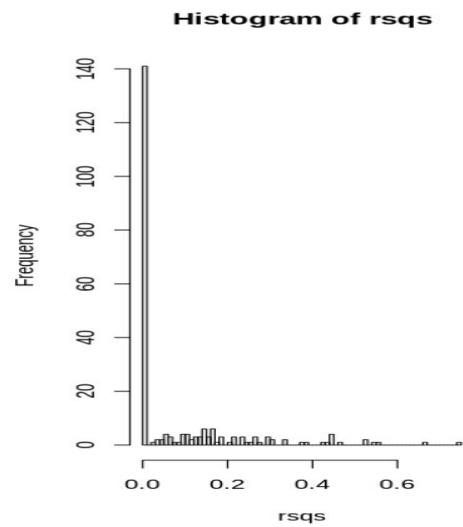
MNA



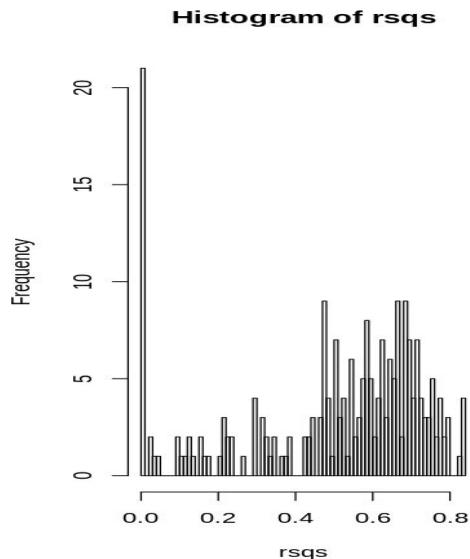
Permutation Test



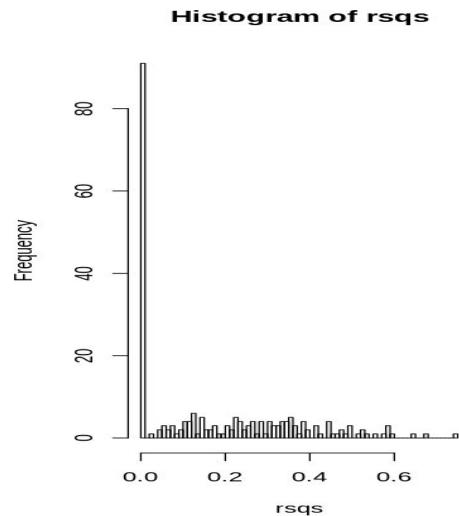
global LASSO $y \sim 13k$ genes



lm + LASSO $y \sim$ selected genes



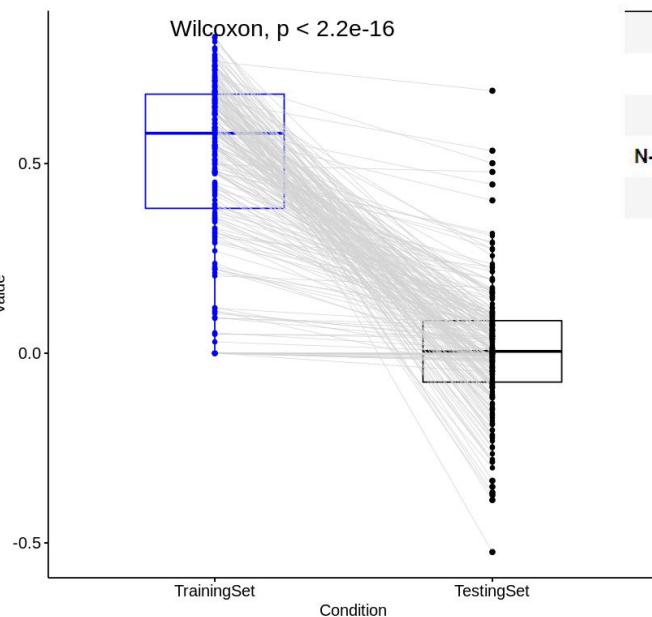
GAM + LASSO $y \sim$ selected genes



LASSO models training



lm + LASSO y~selected genes



R2	group
<dbl>	<chr>
0.769	TrainingSet
0.666	TrainingSet
0.728	TrainingSet
0.497	TrainingSet
0.595	TrainingSet
0.581	TrainingSet
0.745	Training Set
0.484	Training Set
0.497	Training Set
0.455	Training Set
0.497	Training Set
0.387	Training Set

GAM + LASSO y~selected genes

