

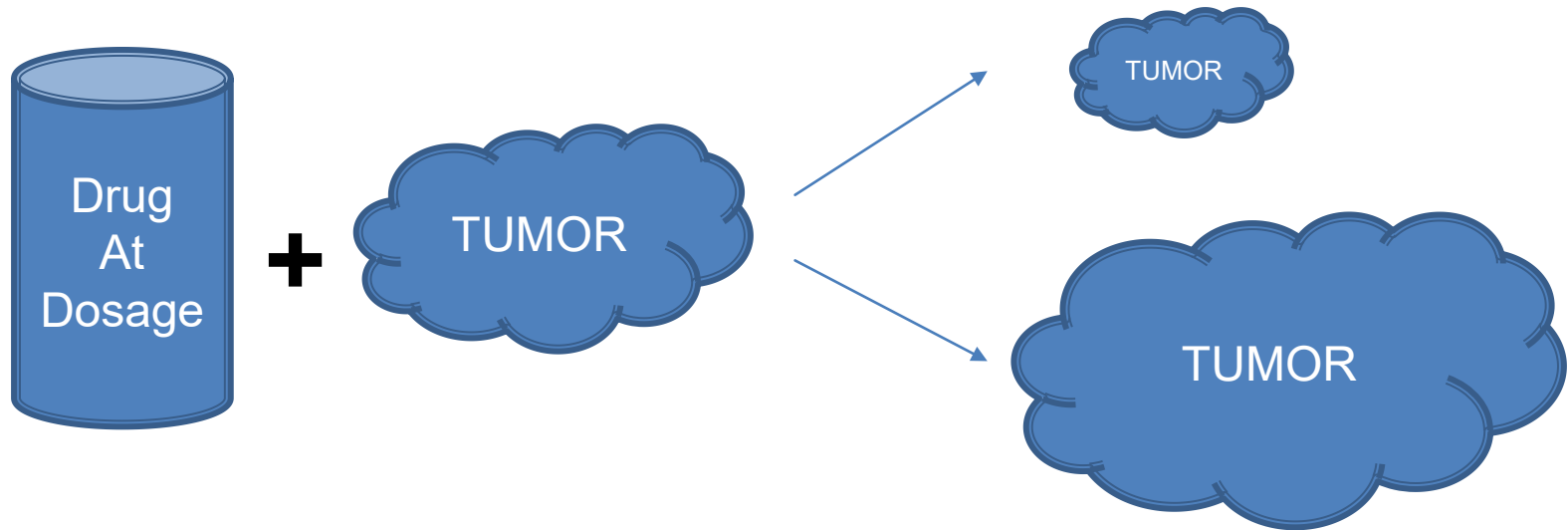
# Semi-Supervised Method for Countering Batch Effect

Stewart He, Jonathan Allen, Ya Ju Fan (LLNL)  
Judith D Cohn (LANL)  
Alex Partin, Fangfang Xia (ANL)



# Prediction Problem

- Given drug features, RNASeq features, and dosage predict tumor reaction.



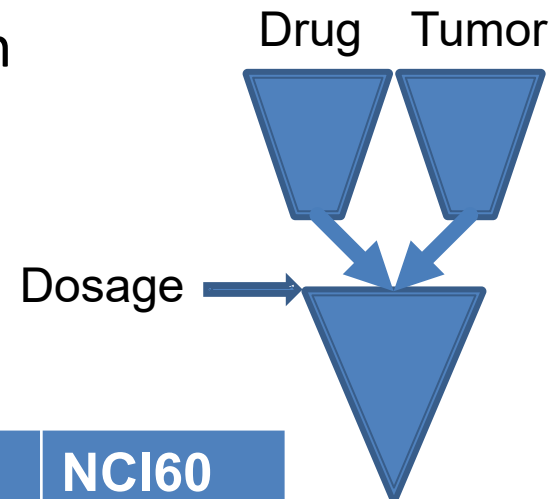
# Data

- Over counts since CCLE and GDSC contain many of the same tumors under different IDs

Dataset	# drugs 3820 dims	# tumors 943 dims	Total datapoints
CCLE	24	474	84,444
CTRP	495	812	4,878,603
gCSI	16	357	49,003
GDSC	239	672	1,100,728
NCI60	1006	59	1,097,284

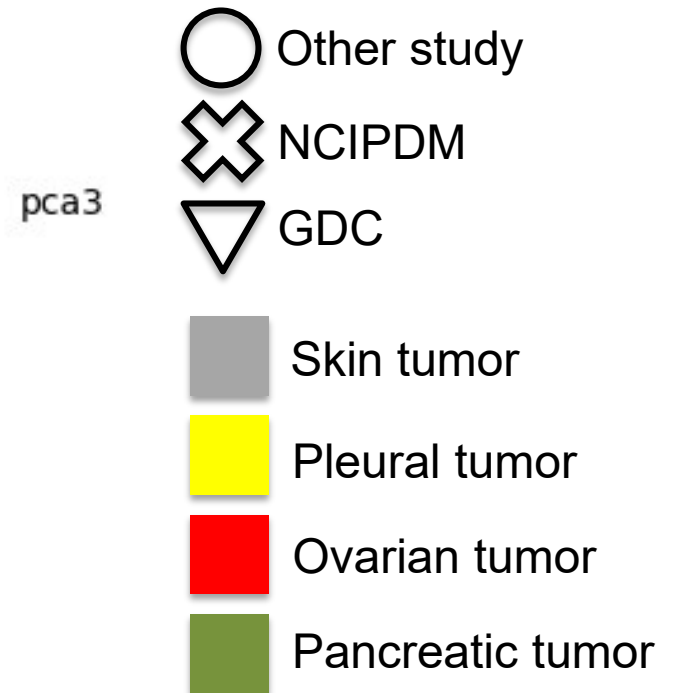
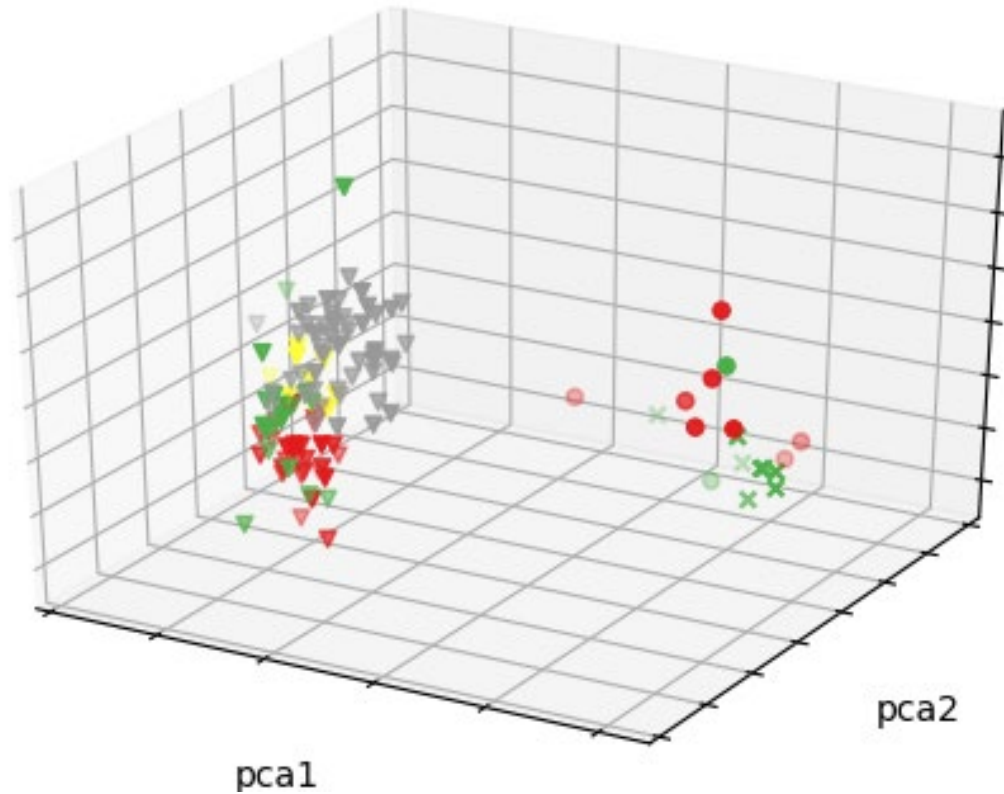
# Regression results

- Used forked neural network to do regression
- Perform inter-dataset test.
  - Example: Train/validate on CCLE test on CTRP
- Poor results using  $R^2$



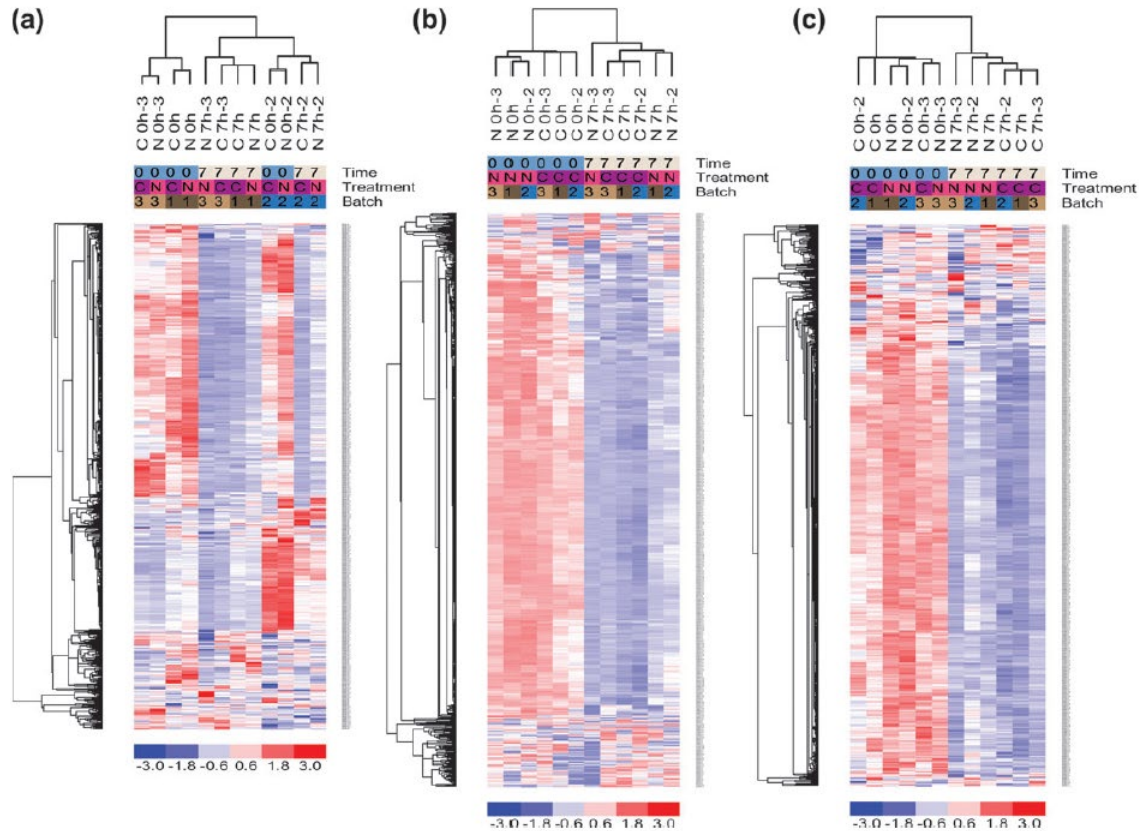
Testing \ Training	CCLE	CTRP	gCSI	GDSC	NCI60
CCLE	0.912	<b>-3.92</b>	<b>-1.28</b>	<b>-3.91</b>	<b>-3.65</b>
CTRP	<b>0.554</b>	0.896	<b>0.454</b>	<b>-.0462</b>	<b>0.06</b>
gCSI	<b>-1.72</b>	<b>-3.40</b>	0.964	<b>-3.54</b>	<b>-2.97</b>
GDSC	<b>-.576</b>	<b>-1.72</b>	<b>0.206</b>	0.879	<b>-1.64</b>
NCI60	<b>0.443</b>	<b>0.226</b>	<b>0.376</b>	<b>0.06</b>	0.894

# PCA across studies



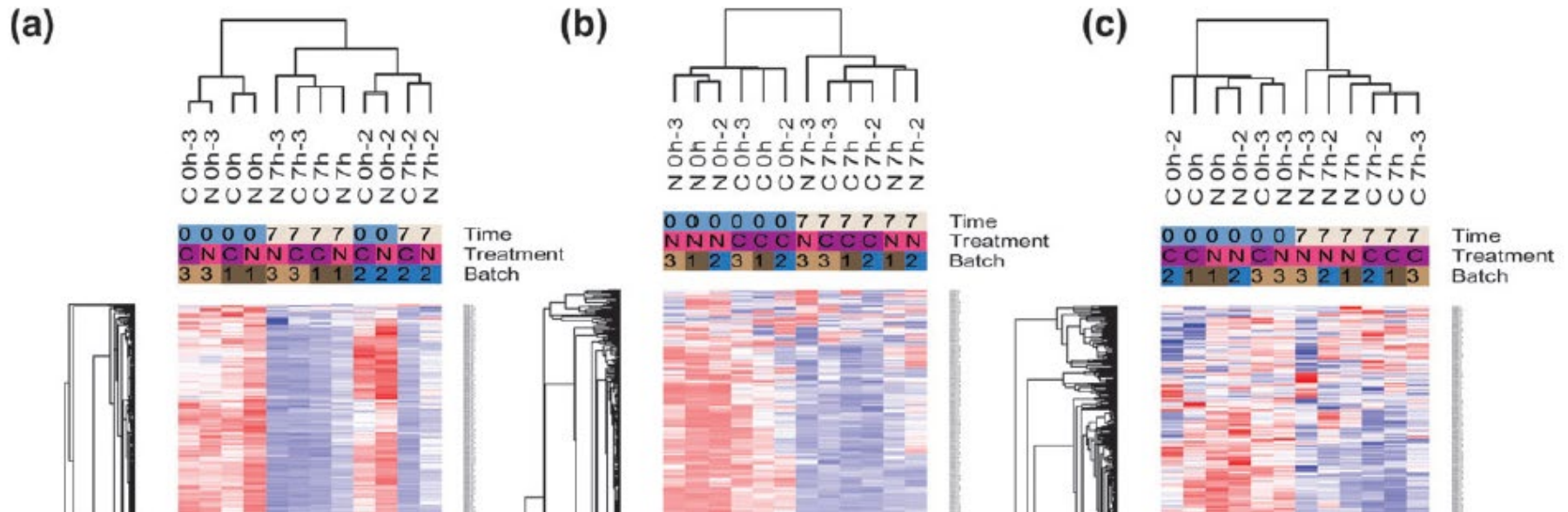
	pca1	pca2	pca3
<b>Explained Variance</b>	<b>0.33</b>	<b>0.12</b>	<b>0.05</b>

# ComBat



Johnson, WE, Rabinovic, A, and Li, C (2007). Adjusting batch effects in microarray expression data using Empirical Bayes methods. *Biostatistics* 8(1):118-127.

# ComBat



Johnson, WE, Rabinovic, A, and Li, C (2007). Adjusting batch effects in microarray expression data using Empirical Bayes methods. *Biostatistics* 8(1):118-127.

# Results: ComBat

Testing Training	CCLE	CTRP	gCSI	GDSC	NCI60
CCLE	0.912	<b>-3.92</b>	<b>-1.28</b>	<b>-3.91</b>	<b>-3.65</b>
CTRP	<b>0.554</b>	0.896	<b>0.454</b>	<b>-0.046</b>	<b>0.06</b>
gCSI	<b>-1.72</b>	<b>-3.40</b>	0.964	<b>-3.54</b>	<b>-2.97</b>
GDSC	<b>-0.576</b>	<b>-1.72</b>	<b>0.206</b>	0.879	<b>-1.64</b>
NCI60	<b>0.443</b>	<b>0.226</b>	<b>0.376</b>	<b>0.06</b>	0.894



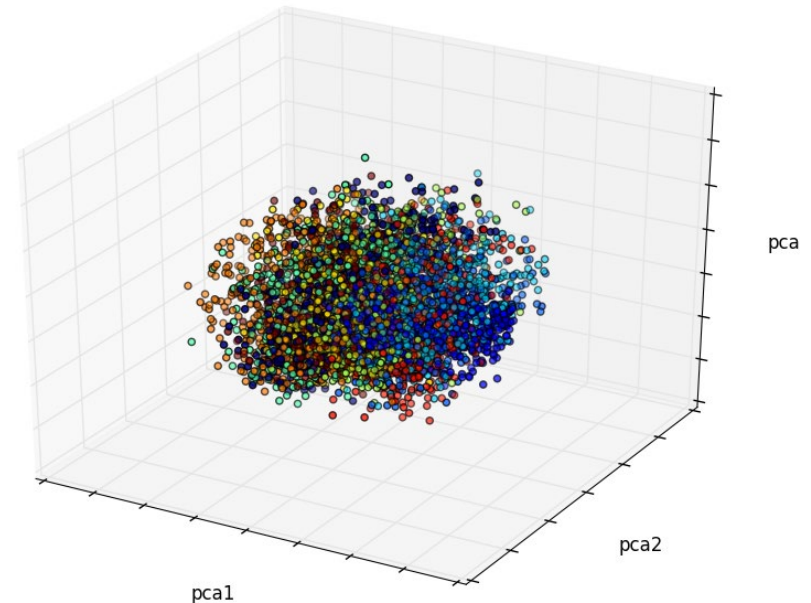
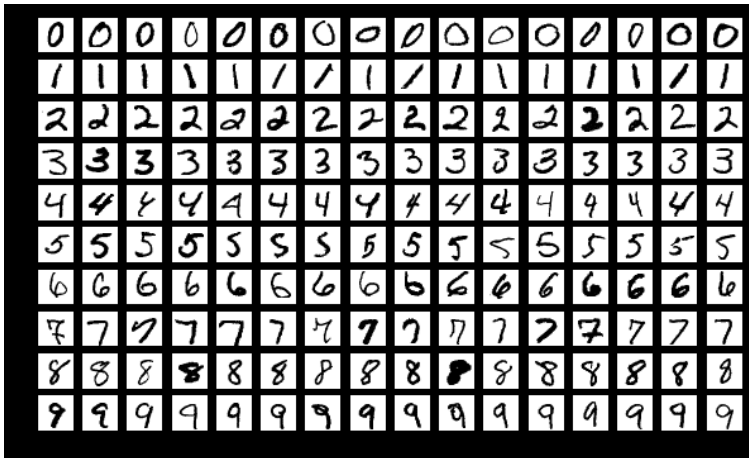
# Learn better features

---

- Very few examples of tumors
- RNASeq originally has 17k features
  - 943 landmark RNASeq genes are hand engineered
- We want better features:
  - More generalized across different types of tumors
  - Learned from unlabeled data
  - Good for regression
  - Use Autoencoders?

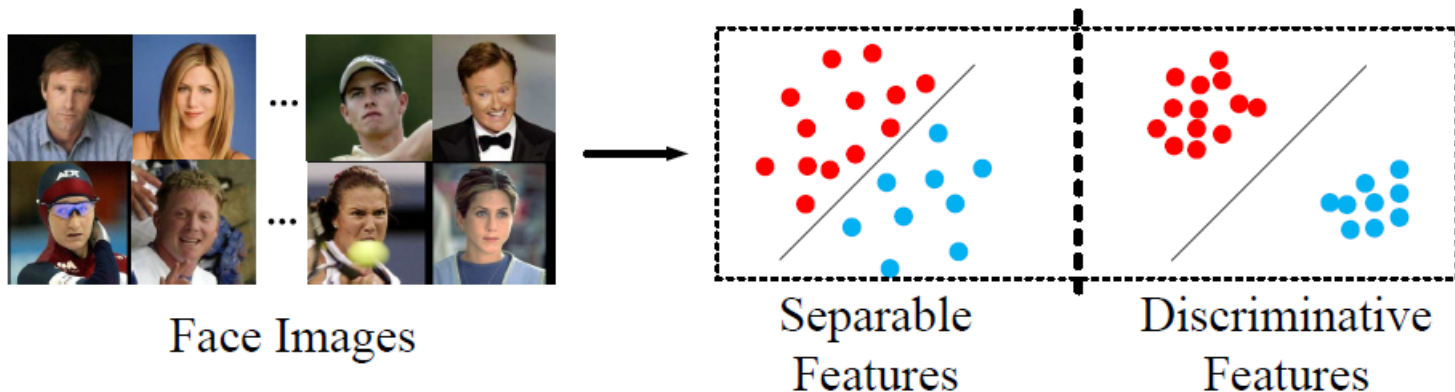
# The problem with Autoencoders

- The latent space is good for reconstruction.
  - That's all the cost function cares about
  - If you're lucky they might be good for other things
- MNIST trained auto encoder
  - First 3 principal components
  - Colored by class



# Modify with Center Loss

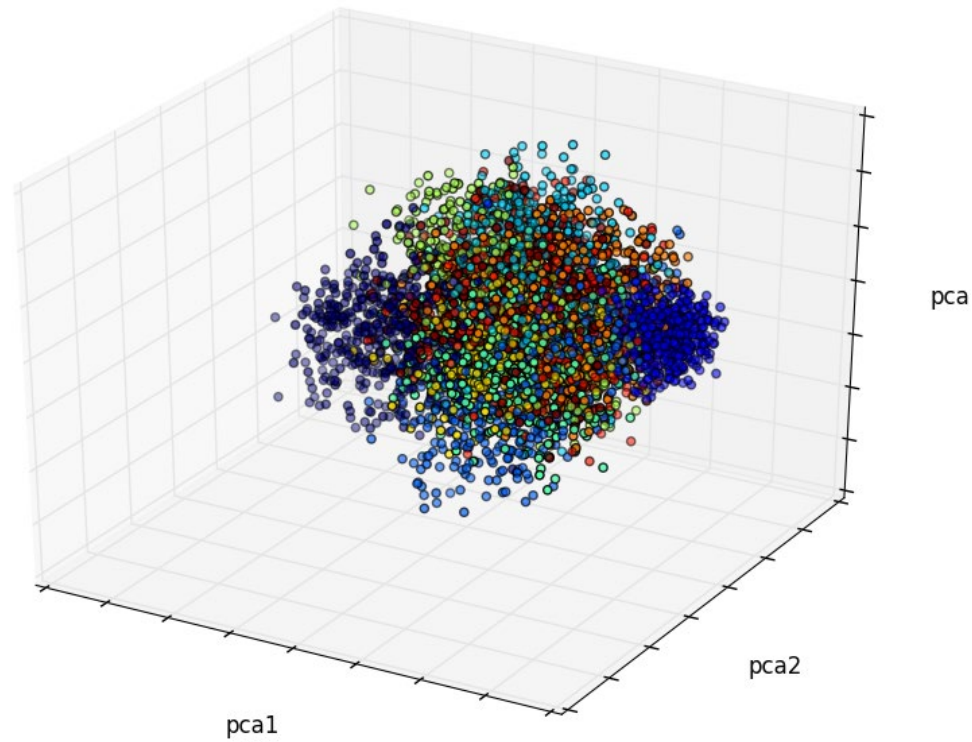
- Center Loss designed to be used for classification
  - A Discriminative Feature Learning Approach for Deep Face Recognition*  
Wen et al., 2016



$$\mathcal{L}_C = \frac{1}{2} \sum_{i=1}^m \|\mathbf{x}_i - \mathbf{c}_{y_i}\|_2^2$$

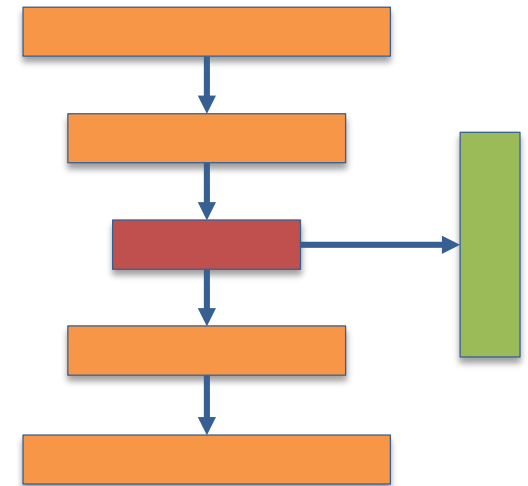
# Center Loss + Autoencoders

- Does not play well. MNIST example:
  - Easily falls for trivial solution



# Classification + Center Loss + Autoencoder

- Autoencoder Bottleneck Classifier
- Must balance 3 terms in loss function
  1. Reconstruction
  2. Classification
  3. Center loss
- Varied the bottleneck between 943 and 20



# Loss comparison between 943 and 20

Feature	Reconstruction Error (RMSE)	F1 score
RNASeq encoded 943	0.716	0.803
RNASeq encoded 20	0.761	0.788
ComBat RNASeq encoded 943	0.714	0.796
ComBat RNASeq encoded 20	0.762	0.788

# Cluster metrics

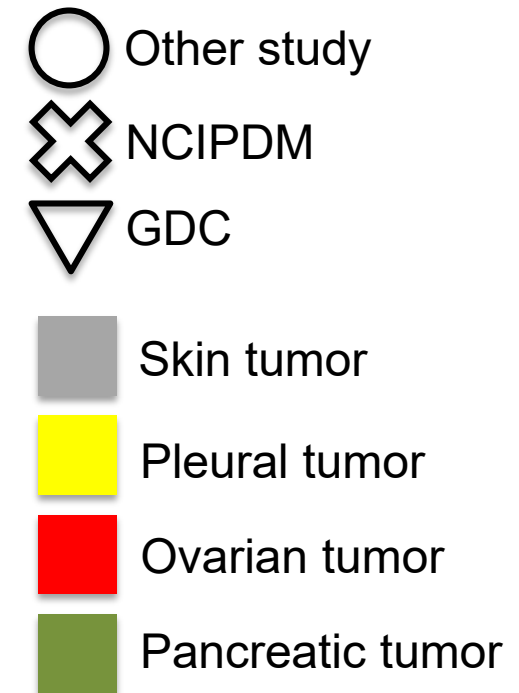
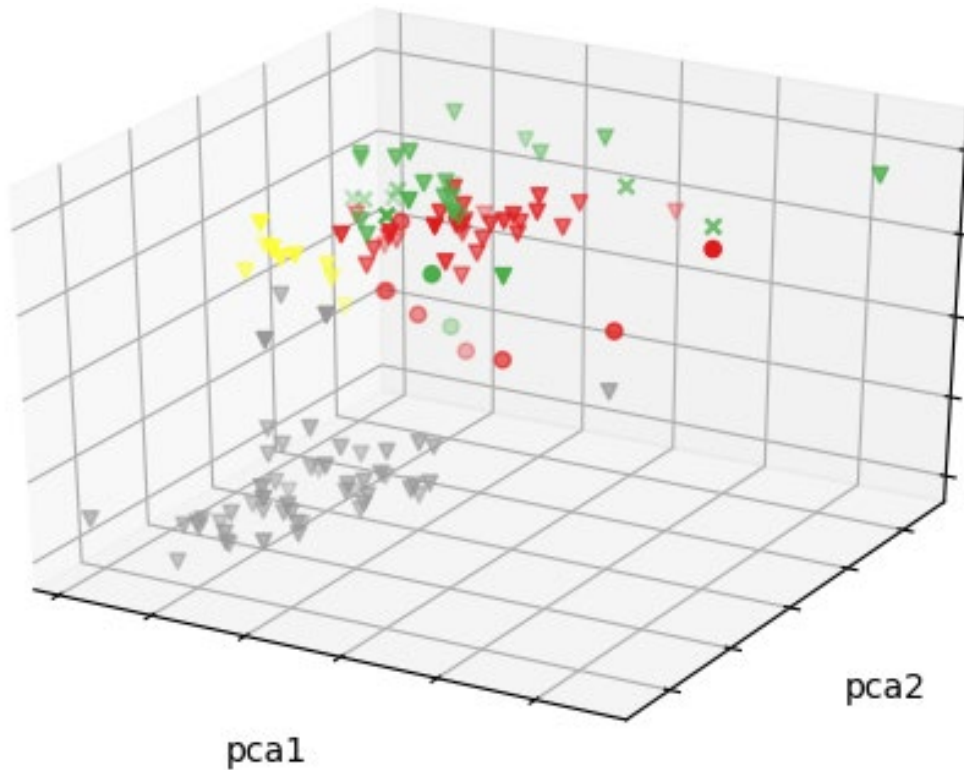
**Silhouette:** between  $[-1, 1]$  where 1 is the best score

**Calinski-Harabasz:** the higher the better

**Davis-Bouldin:** 0 is the best possible score

Dataset	Silhouette	Calinski-Harabasz	Davis-Bouldin
RNASeq	-0.035	16.751	2.558
lincs1000	-0.09	15.012	2.772
ComBat	0.03	20.965	2.51
Encoded RNASeq 943	0.115	53.399	1.881
Encoded RNASeq 20	0.147	110.131	1.689
Encoded ComBat 943	0.105	60.84	1.806
Encoded ComBat 20	<b>0.159</b>	<b>113.451</b>	<b>1.662</b>

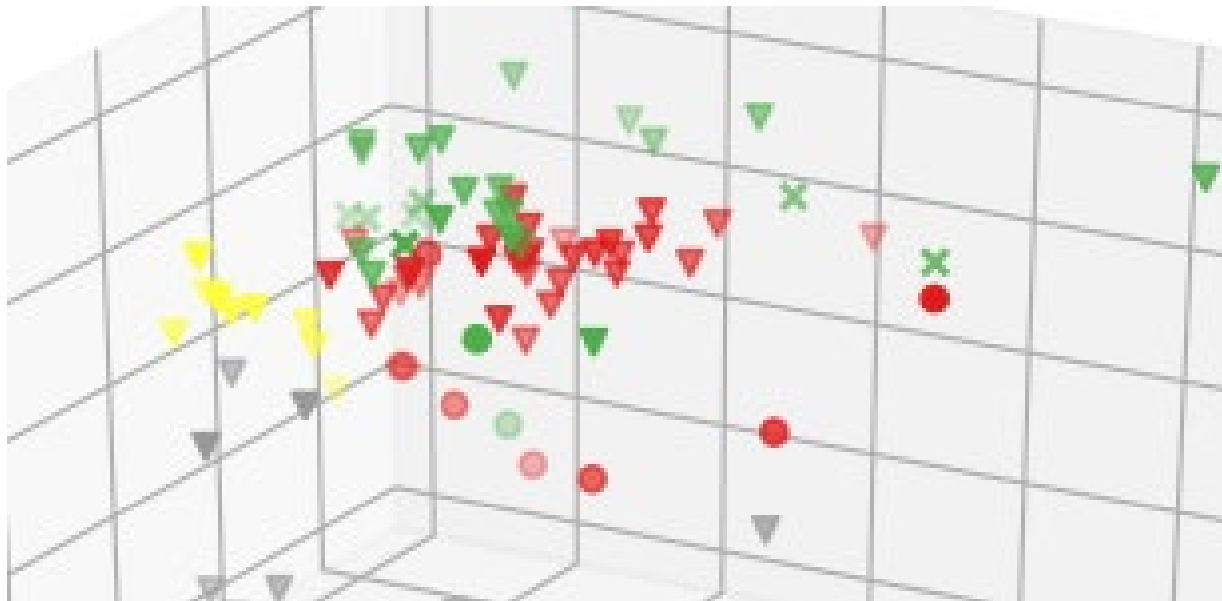
# PCA across studies : AFTER



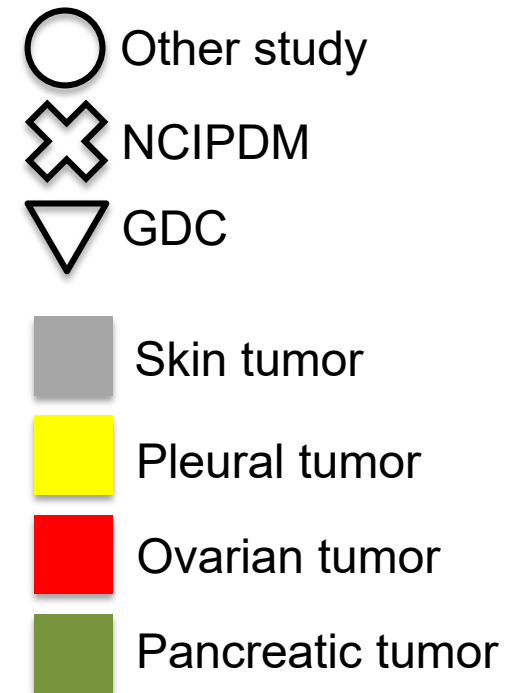
	pca1	pca2	pca3
<b>Explained Variance</b>	<b>0.15</b>	<b>0.11</b>	<b>0.09</b>



# PCA across studies : AFTER



All 3 shapes are now closer together!



	pca1	pca2	pca3
Explained Variance	0.15	0.11	0.09

# Results: Encoded RNASeq 943

Testing \ Training	CCLE	CTRP	gCSI	GDSC	NCI60
CCLE	0.358	<b>0.053</b>	<b>0.304</b>	<b>0.049</b>	<b>0.06</b>
CTRP	<b>0.575</b>	0.681	<b>0.561</b>	<b>0.179</b>	<b>0.366</b>
gCSI	<b>0.048</b>	<b>-0.017</b>	0.206	<b>-0.101</b>	<b>-0.323</b>
GDSC	<b>0.447</b>	<b>-0.049</b>	<b>0.513</b>	0.563	<b>0.13</b>
NCI60	<b>0.435</b>	<b>0.37</b>	<b>0.369</b>	<b>0.23</b>	0.832

# Results: Encoded ComBat 943

Testing \ Training	CCLE	CTRP	gCSI	GDSC	NCI60
CCLE	0.391	<b>-0.025</b>	<b>0.359</b>	<b>0.03</b>	<b>-0.194</b>
CTRP	<b>0.586</b>	0.678	<b>0.566</b>	<b>0.175</b>	<b>0.298</b>
gCSI	<b>0.112</b>	<b>-0.016</b>	0.331	<b>-0.088</b>	<b>-0.489</b>
GDSC	<b>0.419</b>	<b>0.011</b>	<b>0.475</b>	0.549	<b>0.13</b>
NCI60	<b>0.418</b>	<b>0.346</b>	<b>0.388</b>	<b>0.234</b>	0.838

# Results: Encoded RNASeq 20

Testing \ Training	CCLE	CTRP	gCSI	GDSC	NCI60
CCLE	0.727	<b>-0.123</b>	<b>0.359</b>	<b>-0.037</b>	<b>0.044</b>
CTRP	<b>0.599</b>	0.828	<b>0.558</b>	<b>0.132</b>	<b>0.308</b>
gCSI	<b>0.152</b>	<b>0.069</b>	0.37	<b>-0.007</b>	<b>-0.309</b>
GDSC	<b>0.42</b>	<b>-0.104</b>	<b>0.378</b>	0.787	<b>-0.219</b>
NCI60	<b>0.387</b>	<b>0.312</b>	<b>0.405</b>	<b>0.171</b>	0.893

# Results: Encoded ComBat 20

Testing \ Training	CCLE	CTRP	gCSI	GDSC	NCI60
CCLE	0.719	<b>-0.128</b>	<b>0.404</b>	<b>-0.042</b>	<b>-0.051</b>
CTRP	<b>0.574</b>	0.841	<b>0.545</b>	<b>0.082</b>	<b>0.216</b>
gCSI	<b>0.17</b>	<b>0.049</b>	0.376	<b>-0.034</b>	<b>-0.451</b>
GDSC	<b>0.4</b>	<b>-0.072</b>	<b>0.383</b>	0.793	<b>-0.208</b>
NCI60	<b>0.403</b>	<b>0.333</b>	<b>0.377</b>	<b>0.19</b>	0.891

# Extension into PDX data

- Possible use case to train on cell line data and apply to PDX data.
  - Spearman rank correlation to relate PDX doubling times and AUC dosage/response rate curve

	ComBat	Encoded RNASeq 943
CCLE	-0.07	-0.07
CTRP	-0.24	-0.304
gCSI	-0.23	-0.261
GDSC	0	-0.033
NCI60	-0.075	-0.077

# Acknowledgements

---

- Maulik Shukla - data organization/access
- Ben McMahon - helped Judith with clustering
- Rick Stevens - project PI
- Austin Clyde – performing regression tests
- Jason David Gans – doubling time predictions on PDX dataset





# Learned features applied to regression

- Modest improvement to regression results

Testing-> Trianing	CCLE	CTRP	gCSI	GDSC	NCI60
CCLE	.913	<b>-1.394</b>	<b>-.061</b>	<b>-1.477</b>	<b>-.640</b>
CTRP	<b>.611</b>	.817	<b>.501</b>	<b>.136</b>	<b>.186</b>
gCSI	<b>.021</b>	<b>-.357</b>	.962	<b>-.201</b>	<b>-.906</b>
GDSC	<b>.380</b>	<b>-.445</b>	<b>.364</b>	.861	<b>-.238</b>
NCI60	<b>.450</b>	<b>.270</b>	<b>.351</b>	<b>.160</b>	.894

# Supplemental full image

