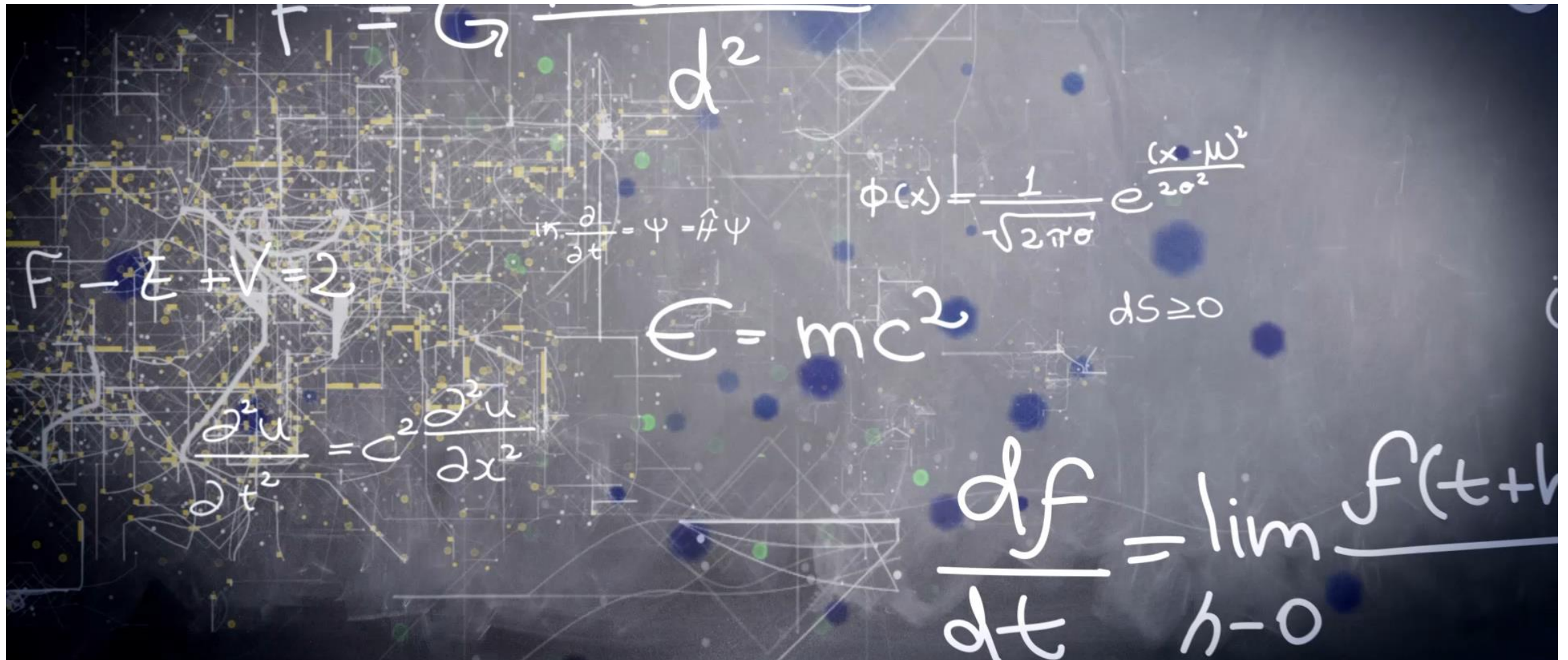


# Machine Learning



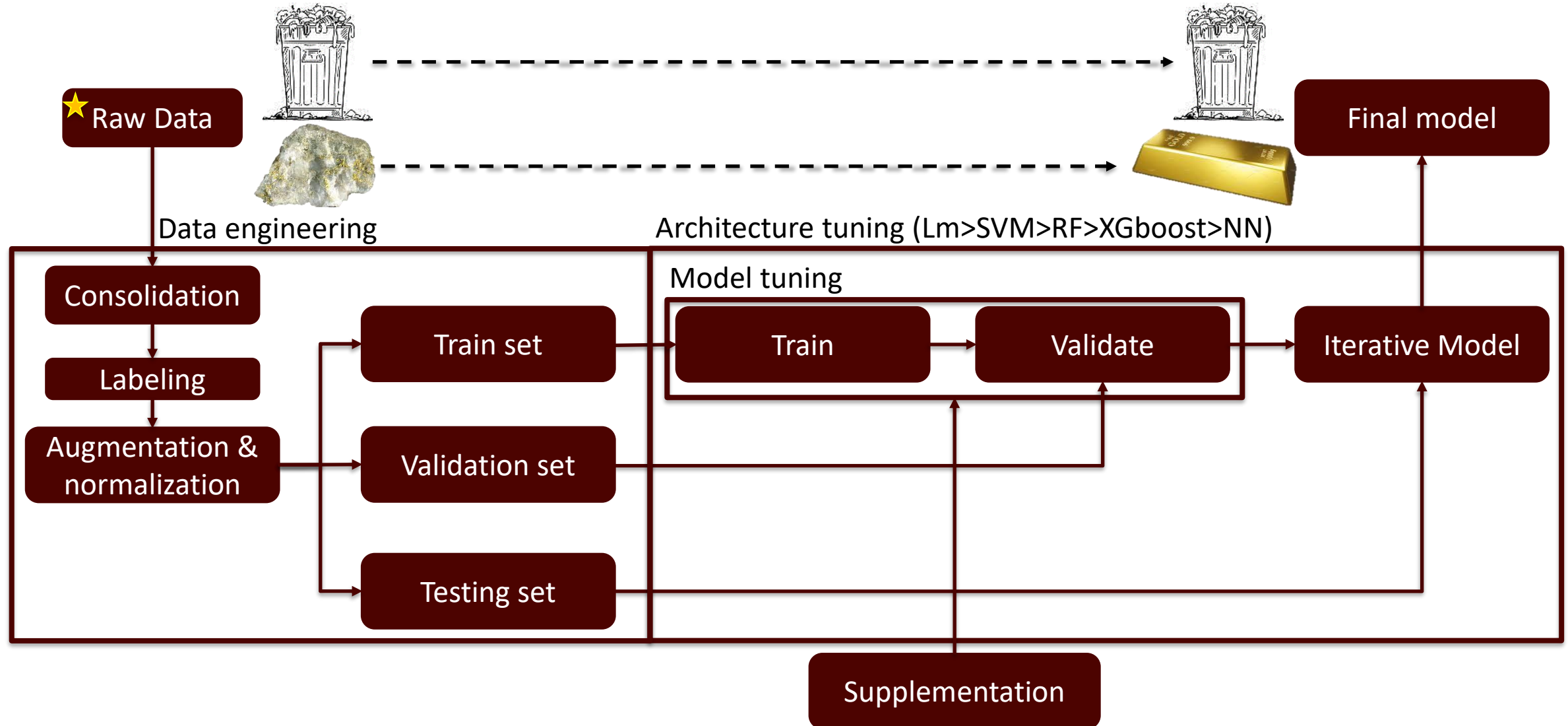
TEXAS A&M  
UNIVERSITY



# Training a supervised model



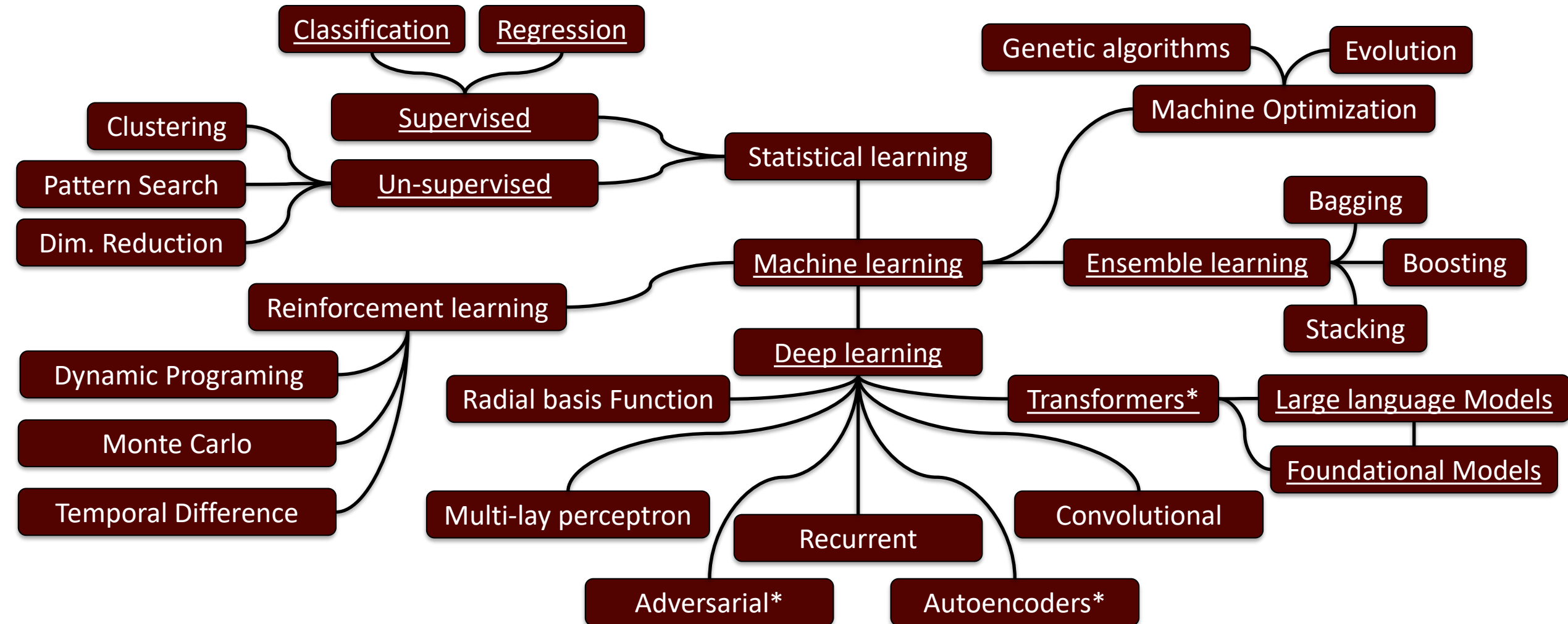
TEXAS A&M  
UNIVERSITY



# A network of ML methods



TEXAS A&M  
UNIVERSITY



Underline indicts linked information in the appendix. \* indicates generative AI



The conventional bar to entry...



TEXAS A&M  
UNIVERSITY



# Low-code AutoML (H2O Flow)



TEXAS A&M  
UNIVERSITY

```
Python 3.8.17 (default, Jul 5 2023, 20:44:21) [MSC v.1916 64 bit (AMD64)] :: Anaconda, Inc. on win32
Type "help", "copyright", "credits" or "license" for more information.
>>> import h2o
>>> h2o.init()
```

H<sub>2</sub>O FLOW Flow- Cell- Data- Model- Score- Admin- Help-

Untitled Flow



assist

## ? Assistance

Routine	Description
<a href="#">importFiles</a>	Import file(s) into H <sub>2</sub> O
<a href="#">importSqlTable</a>	Import SQL table into H <sub>2</sub> O
<a href="#">getFrames</a>	Get a list of frames in H <sub>2</sub> O
<a href="#">splitFrame</a>	Split a frame into two or more frames
<a href="#">mergeFrames</a>	Merge two frames into one
<a href="#">getModels</a>	Get a list of models in H <sub>2</sub> O
<a href="#">getGrids</a>	Get a list of grid search results in H <sub>2</sub> O
<a href="#">getPredictions</a>	Get a list of predictions in H <sub>2</sub> O
<a href="#">getJobs</a>	Get a list of jobs running in H <sub>2</sub> O
<a href="#">runAutoML</a>	Automatically train and tune many models
<a href="#">buildModel</a>	Build a model
<a href="#">importModel</a>	Import a saved model
<a href="#">predict</a>	Make a prediction

CS

importFiles

## Import Files

Search: Enter a file or directory path and press the Enter key

Selected Files: (No files selected)

Actions: Import

OUTLINE FLOWS CLIPS **HELP**

## Help

Using Flow for the first time?

Quickstart Videos

Or, [view example Flows](#) to explore and learn H<sub>2</sub>O.

STAR H2O ON GITHUB!

Star

GENERAL

- [Flow Web UI ...](#)
- [... Importing Data](#)
- [... Building Models](#)
- [... Making Predictions](#)
- [... Using Flows](#)
- [... Troubleshooting Flow](#)

EXAMPLES

Flow packs are a great way to explore and learn H<sub>2</sub>O. Try out these Flows and run them in your browser.  
[Browse installed packs...](#)

H<sub>2</sub>O REST API

- [Routes](#)
- [Schemas](#)

# Unsupervised ML

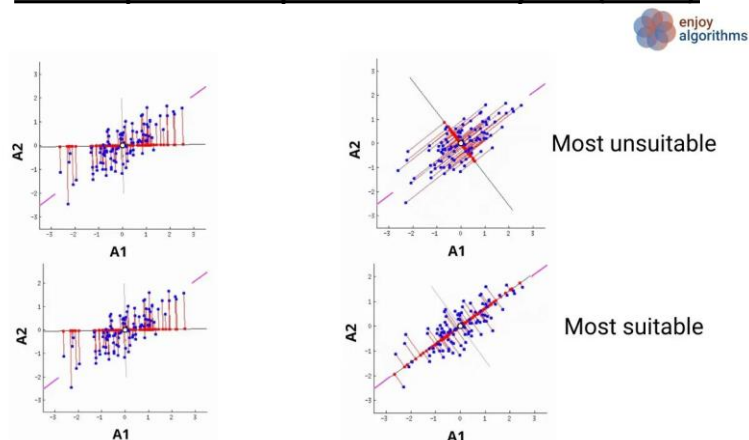


TEXAS A&M  
UNIVERSITY

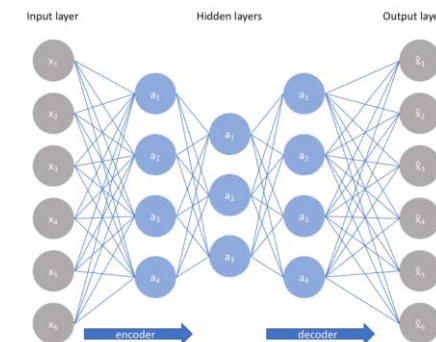
## Hierarchical Clustering



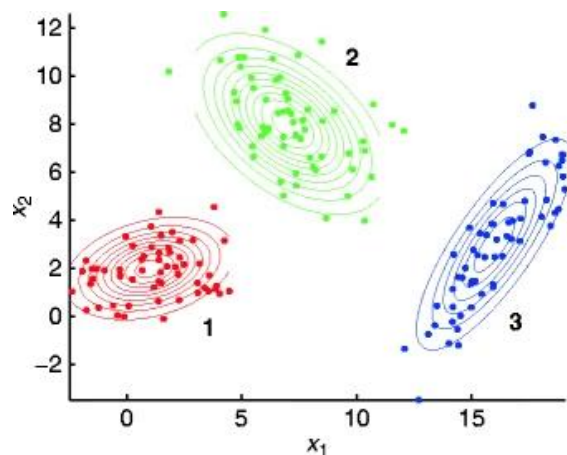
## Principle component analysis (PCA)



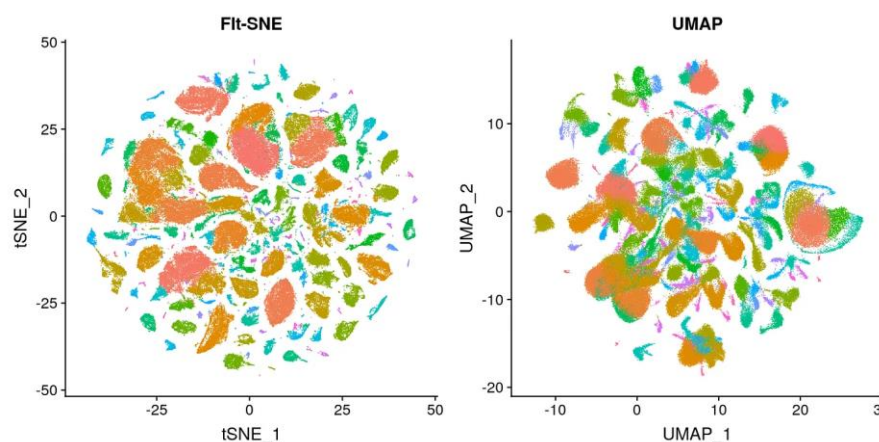
## Autoencoders



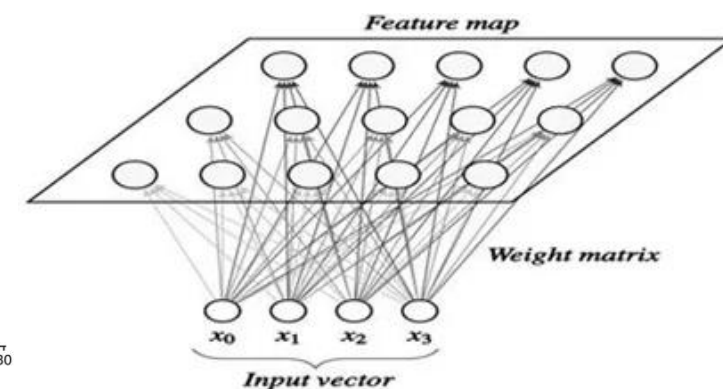
## Model-based clustering



## t-SNE & UMAP



## Self organizing maps (SOM)

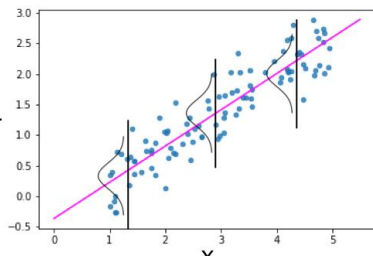


\*\*Citations/info links in images\*\*

# Conventional supervised ML



GLM



Lasso/Ridge/Elastic network

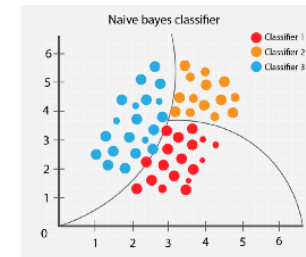
$$\frac{\sum_{i=1}^n (y_i - x_i^T \hat{\beta})^2}{2n} + \lambda \left( \frac{1-\alpha}{2} \sum_{j=1}^m \hat{\beta}_j^2 + \alpha \sum_{j=1}^m |\hat{\beta}_j| \right)$$

Naïve Bayes model

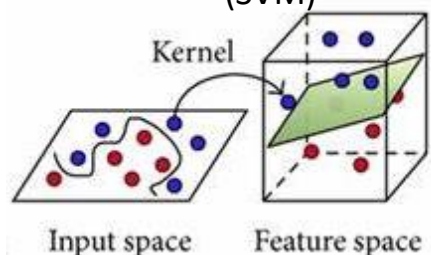
$$P(A|B) = \frac{P(B|A) P(A)}{P(B)}$$

using Bayesian probability terminology, the above equation can be written as

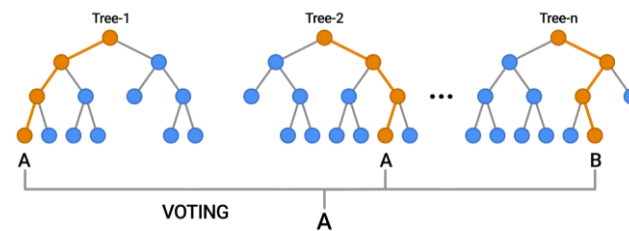
$$\text{Posterior} = \frac{\text{prior} \times \text{likelihood}}{\text{evidence}}$$



Support vector machine (SVM)



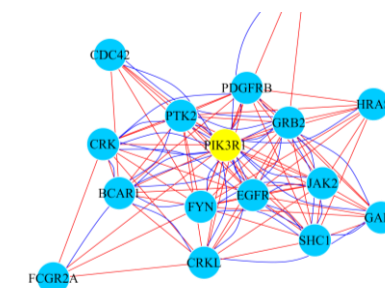
Tree/RF/XGBoost



Recommender systems



Network-based algorithms





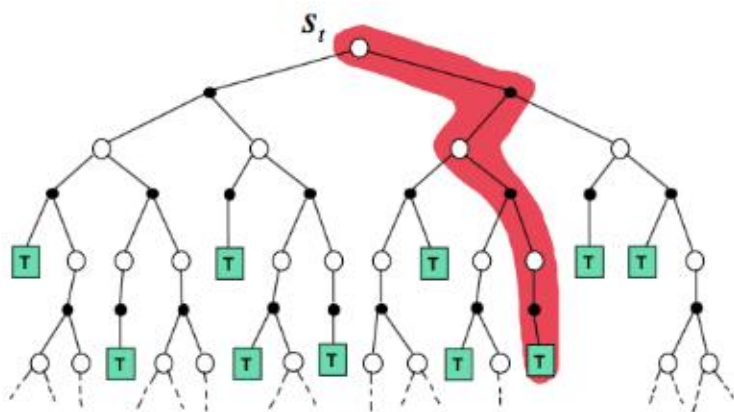
# Reinforcement learning



TEXAS A&M  
UNIVERSITY

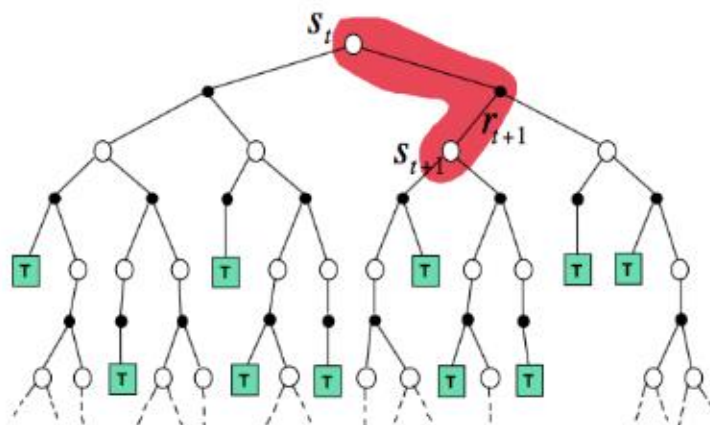
Monte-Carlo

$$V(S_t) \leftarrow V(S_t) + \alpha (G_t - V(S_t))$$



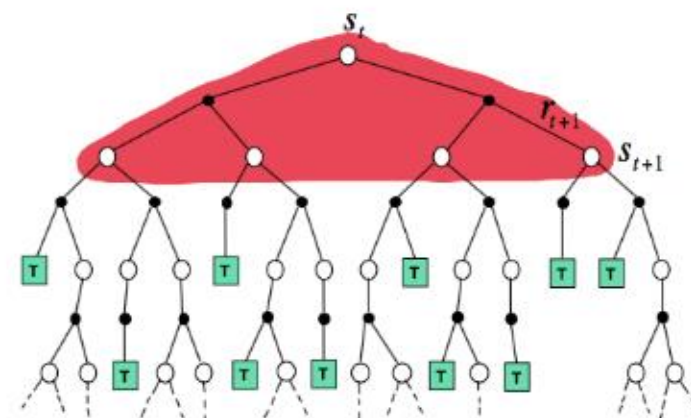
Temporal-Difference

$$V(S_t) \leftarrow V(S_t) + \alpha (R_{t+1} + \gamma V(S_{t+1}) - V(S_t))$$



Dynamic Programming

$$V(S_t) \leftarrow \mathbb{E}_\pi [R_{t+1} + \gamma V(S_{t+1})]$$





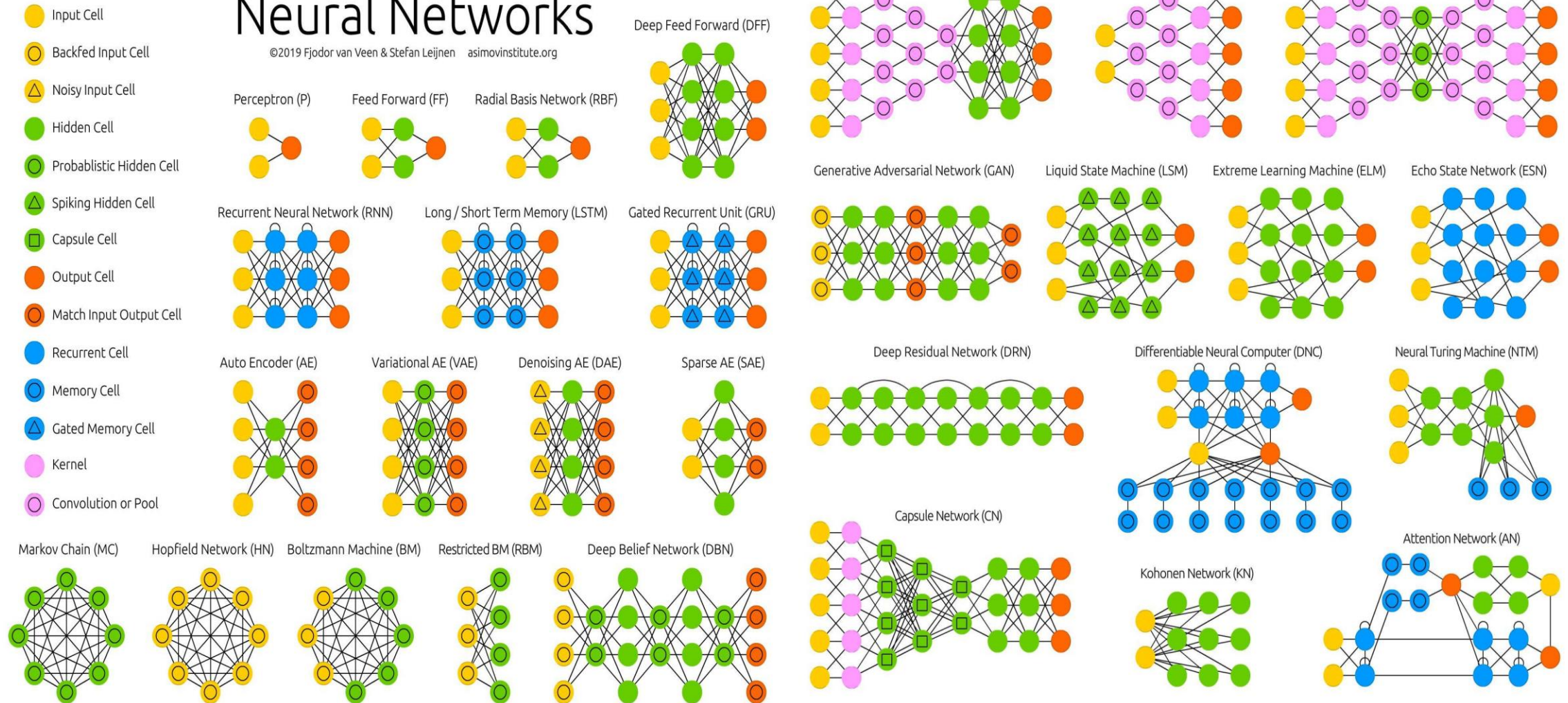
# Deep learning models



TEXAS A&M  
UNIVERSITY

## A mostly complete chart of Neural Networks

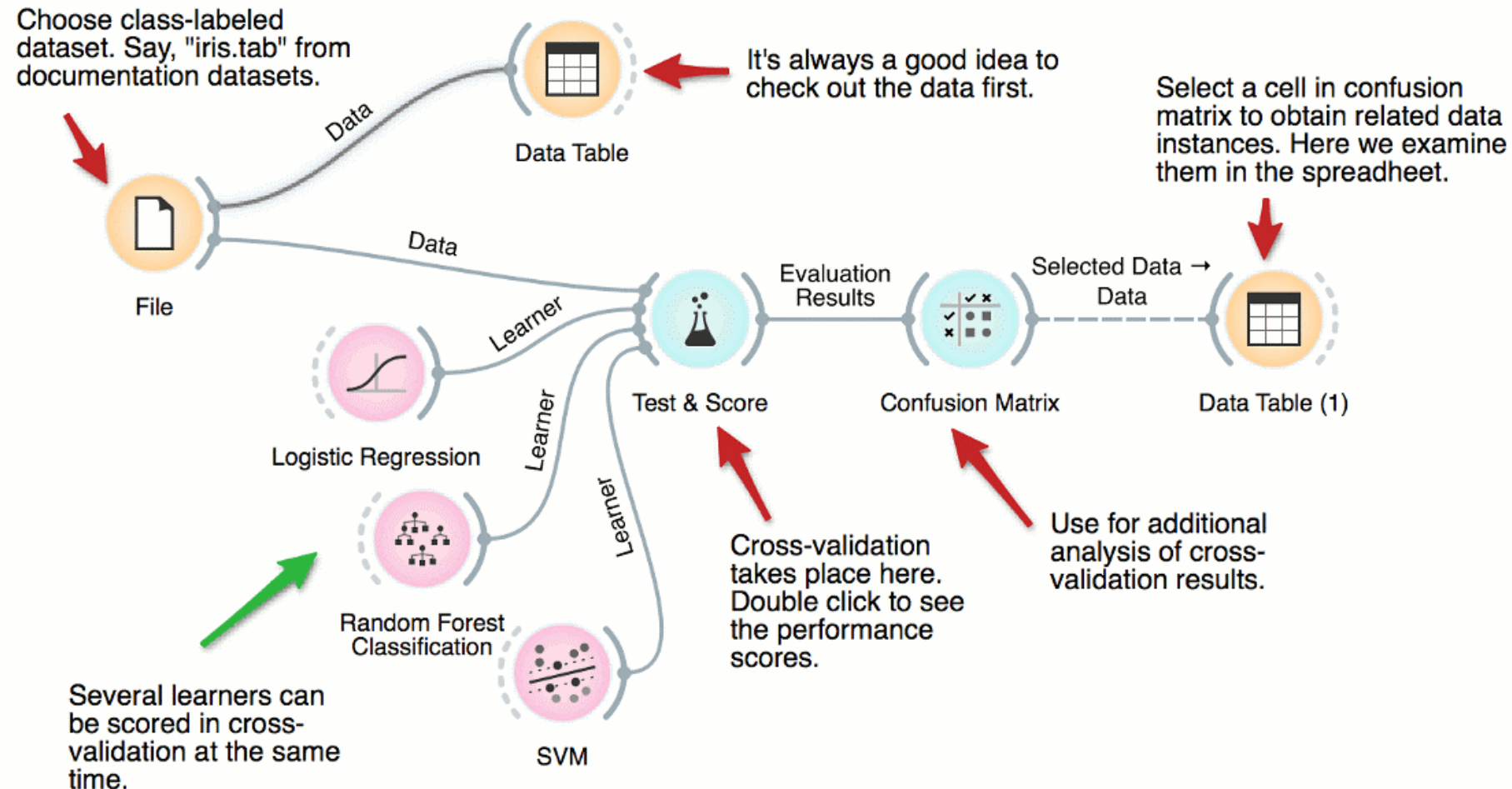
©2019 Fjodor van Veen & Stefan Leijnen asimovinstitute.org



# No-Code Generalized ML



TEXAS A&M  
UNIVERSITY



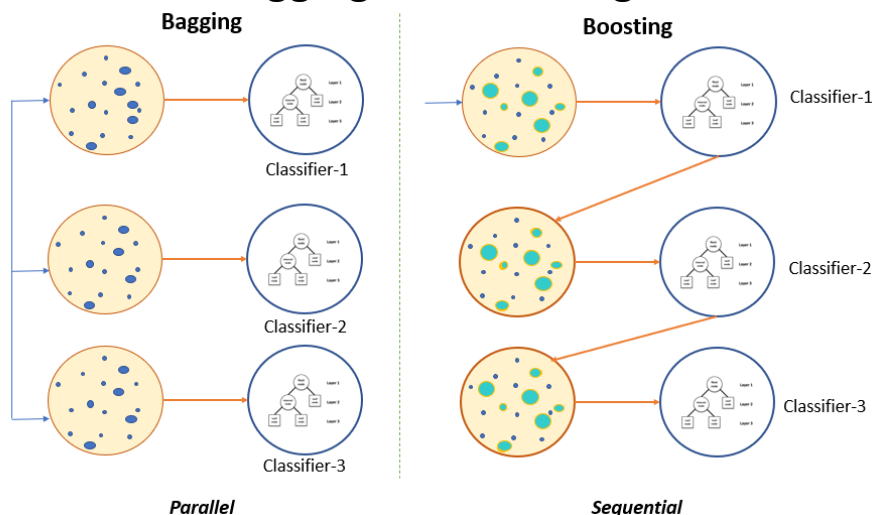


# Ensemble modeling



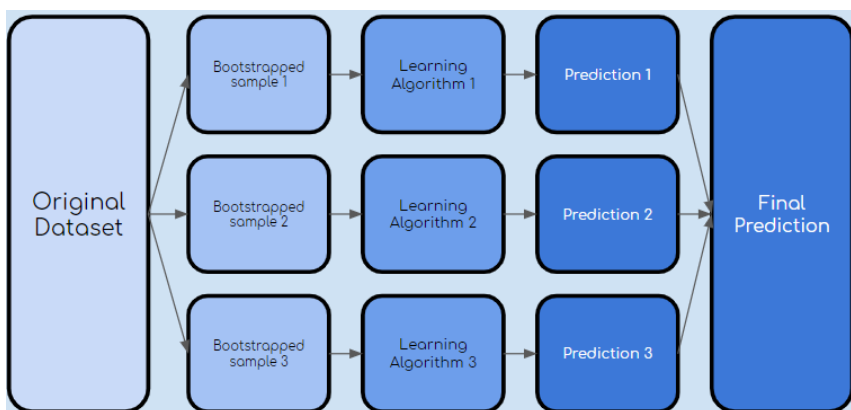
TEXAS A&M  
UNIVERSITY

## Bagging and boosting



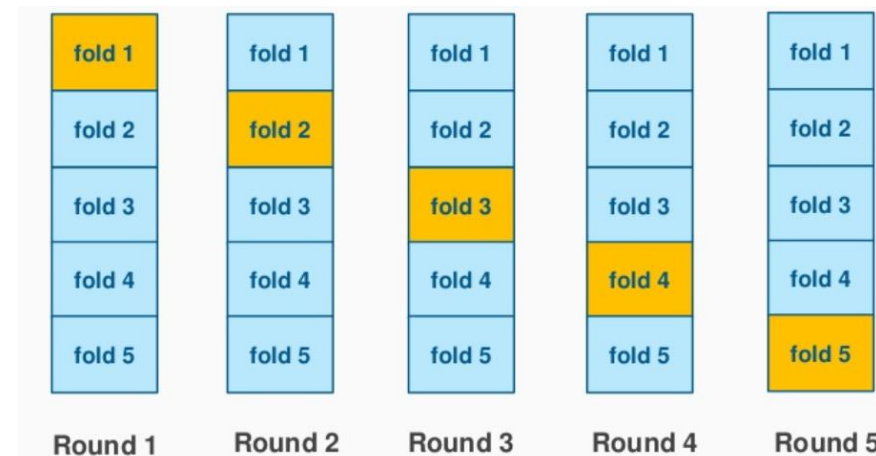
Performance enhancing

## Stacking

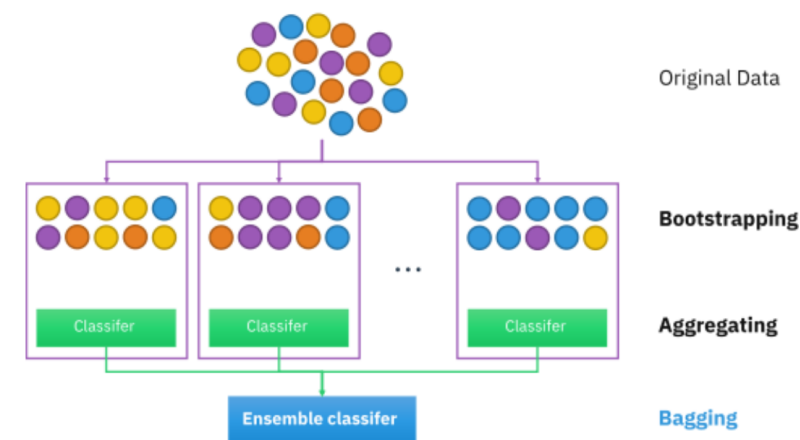


Robustness enhancing

## Cross-validated



## Bootstrapped

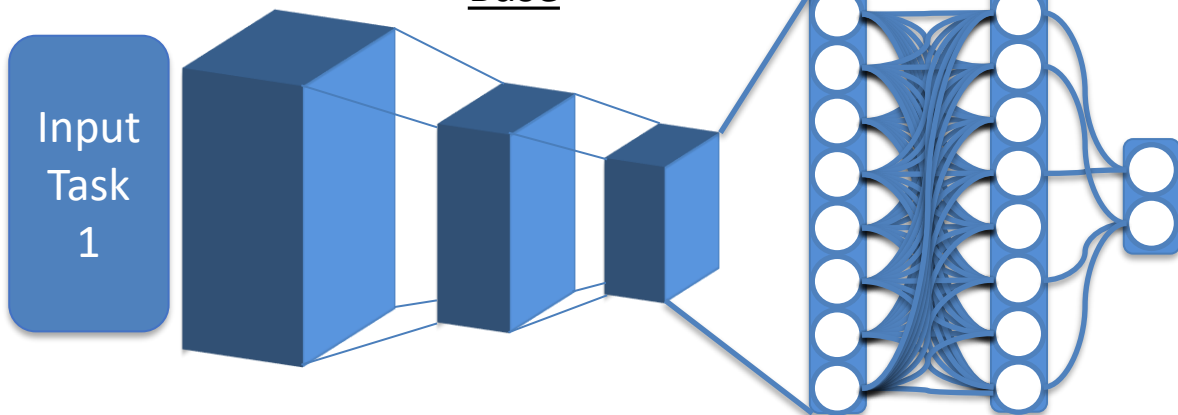


# Transfer learning methods

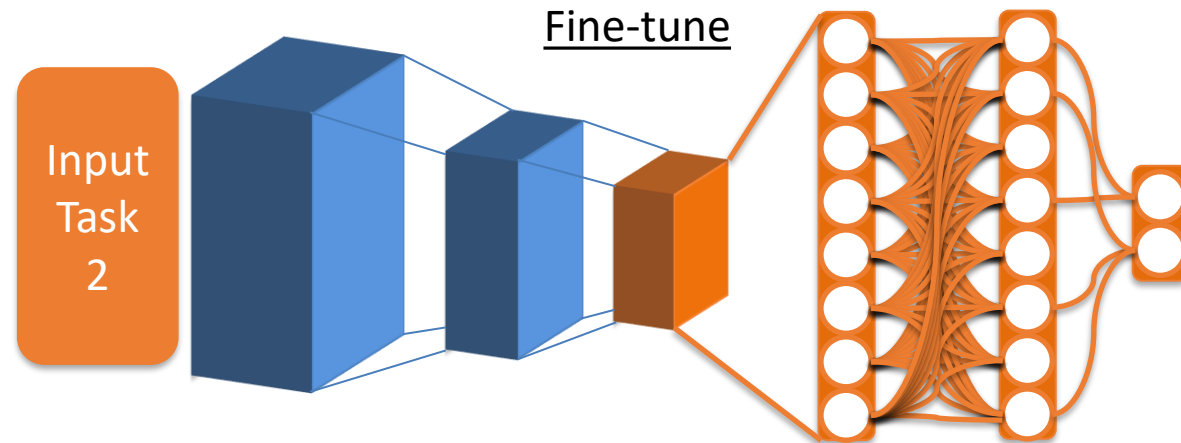


TEXAS A&M  
UNIVERSITY

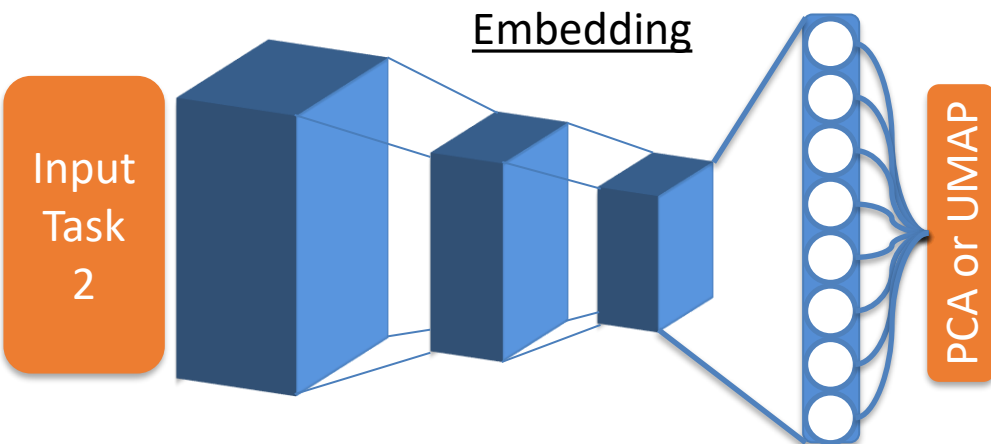
Base



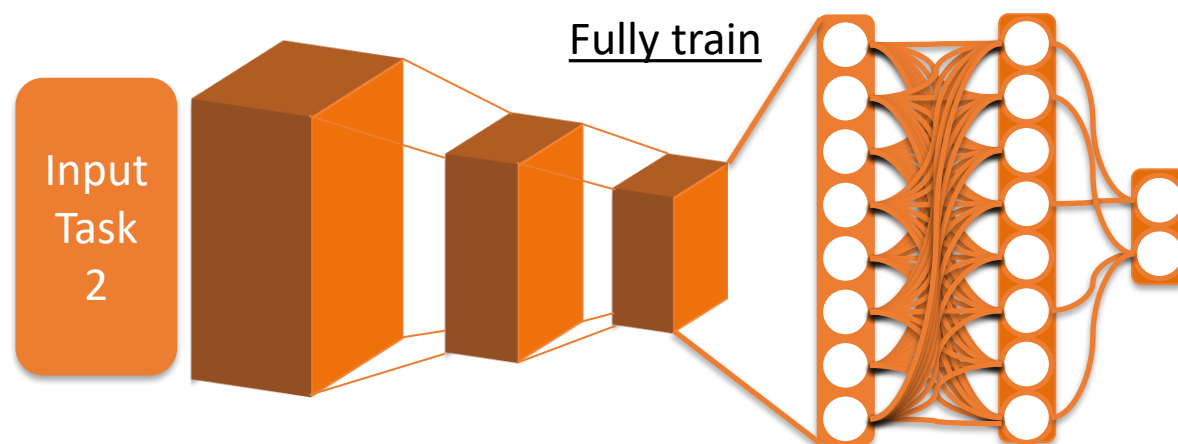
Fine-tune



Embedding

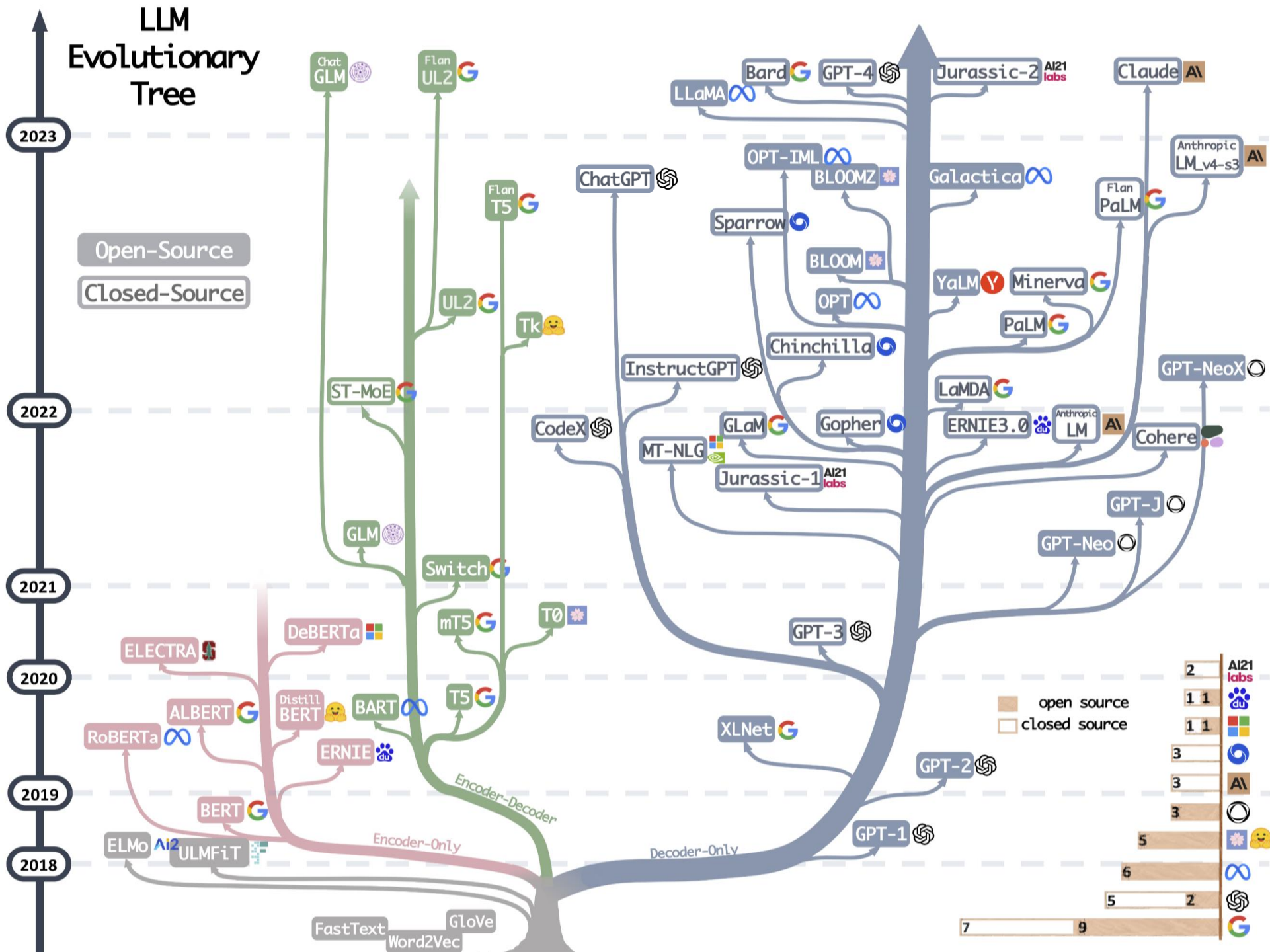


Fully train





# LLM Evolutionary Tree



2	AI21 labs
1 1	du
1 1	du
3	
3	AI
3	
5	
6	
5	2
7	9

# Foundational models



TEXAS A&M  
UNIVERSITY

**Broad Training Data**: trained on extensive datasets, which require substantial computational resources. This training allows them to learn a wide range of tasks and skills during the initial phase.



**Self-Supervision**: Generally, use self-supervision techniques during training where labels or targets are generated from the data itself, rather than relying solely on human-labeled data.



**Large Parameter Count**: Typically contains at least billions of parameters to enable them to capture complex patterns and relationships in the data.



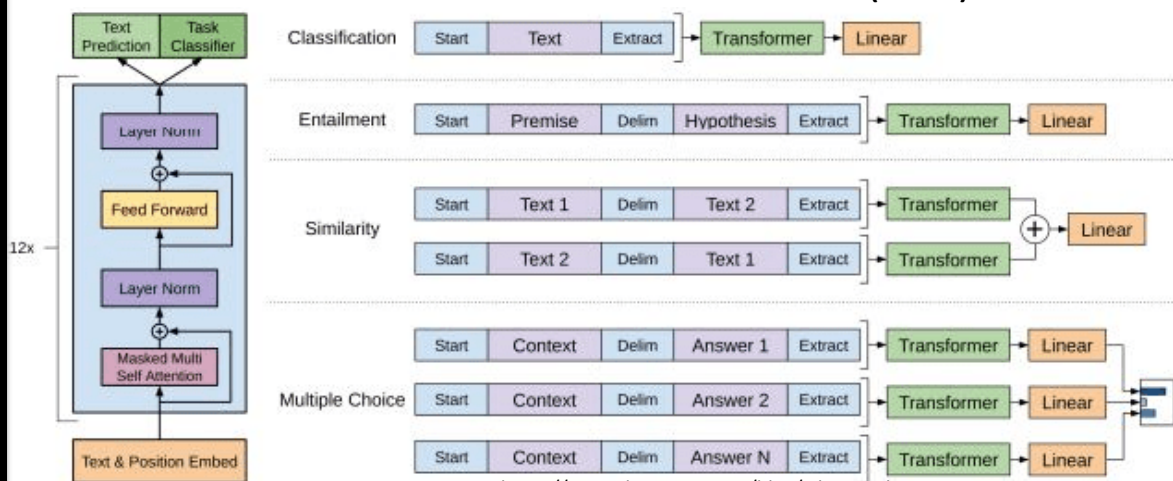
**Applicability Across Contexts**: Applicable across a wide range of contexts, can be secondarily fine-tuned for specific tasks with minimal adjustments, making them highly versatile.

# Foundation models



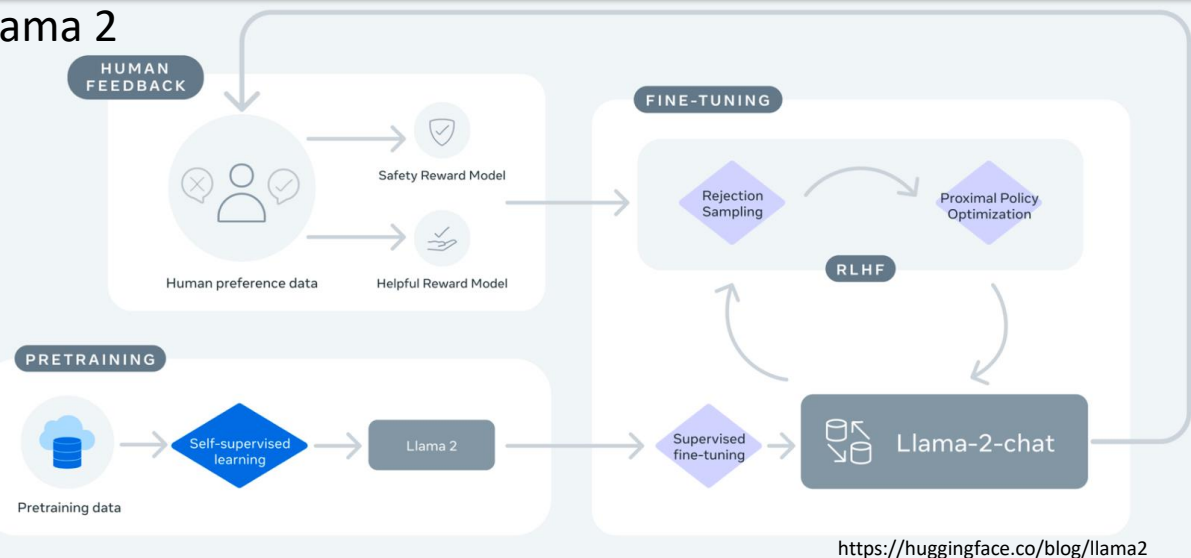
TEXAS A&M  
UNIVERSITY

## Generative Pre-trained Transformer (GPT)



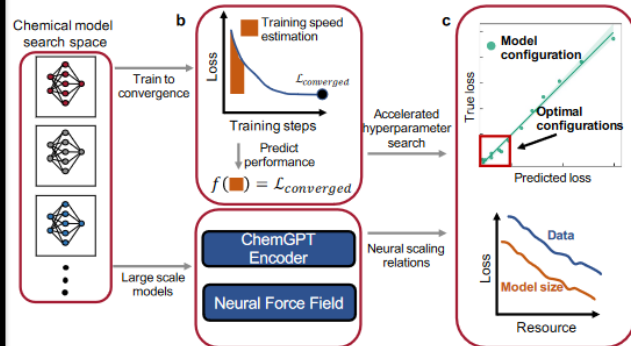
<https://www.datacamp.com/blog/what-we-know-gpt4>

## Llama 2



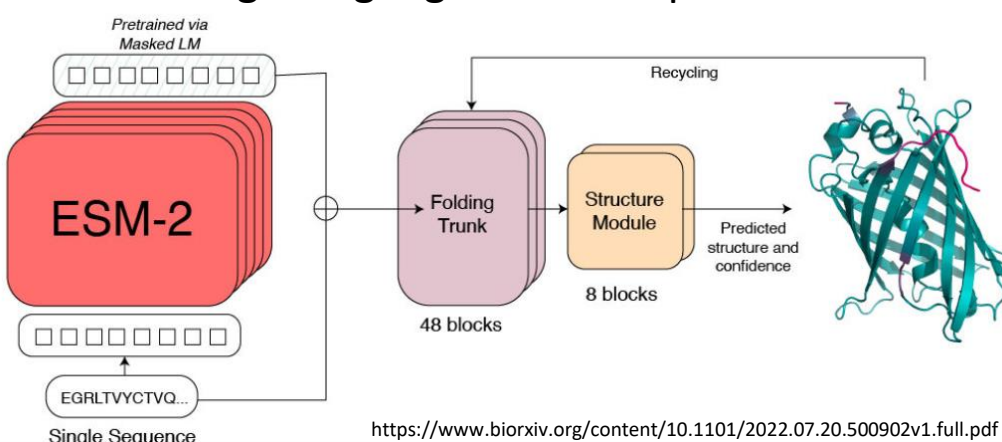
<https://huggingface.co/blog/llama2>

## ChemGPT



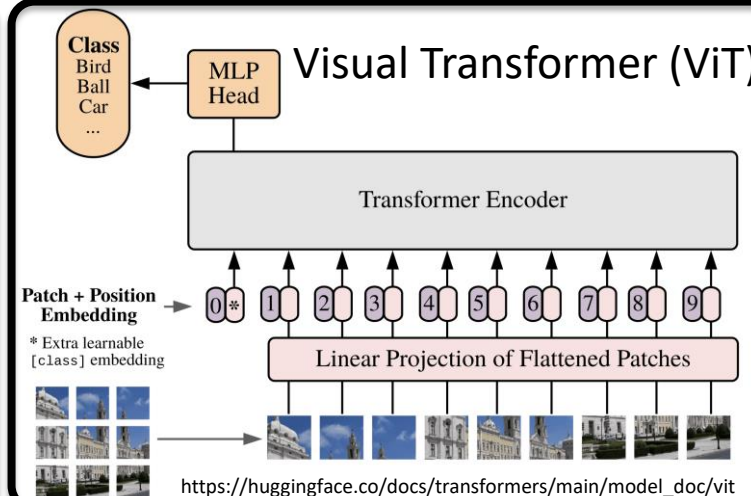
<https://chemrxiv.org/engage/chemrxiv/article-details/627bdd544bdd532395fb4b5>

## ESM2: A large language model of protein structure



<https://www.biorxiv.org/content/10.1101/2022.07.20.500902v1.full.pdf>

## Visual Transformer (ViT)



[https://huggingface.co/docs/transformers/main/model\\_doc/vit](https://huggingface.co/docs/transformers/main/model_doc/vit)

Open source and local deployment demo's: Hugging Face



Model Type	Implementation & Key Features	Use Cases
Memory Networks	<ul style="list-style-type: none"><li>- LSTM</li><li>- External memory for handling long-term dependencies</li><li>- Read and write operations.</li></ul>	<ul style="list-style-type: none"><li>- Question-answering</li><li>- Dialogue systems</li></ul>
Causal Language Modeling (CLM)	<ul style="list-style-type: none"><li>- GPT, Llama</li><li>- Autoregressive model that predicts sequential tokens.</li><li>- Unidirectional context.</li></ul>	<ul style="list-style-type: none"><li>- Text generation</li><li>- Summarization</li></ul>
Masked Language Modeling (MLM)	<ul style="list-style-type: none"><li>- BERT, RoBERTa</li><li>- Input tokens are masked</li><li>- Model predicts context.</li><li>- Bidirectional context.</li></ul>	<ul style="list-style-type: none"><li>- Text classification</li><li>- Sentiment analysis</li><li>- Named entity recognition</li></ul>
Sequence-to-Sequence (Seq2Seq)	<ul style="list-style-type: none"><li>- T5</li><li>- Encoder-decoder architecture.</li><li>- Handles input-output transformations.</li></ul>	<ul style="list-style-type: none"><li>- Machine translation</li><li>- Summarization</li><li>- Question-answering</li></ul>



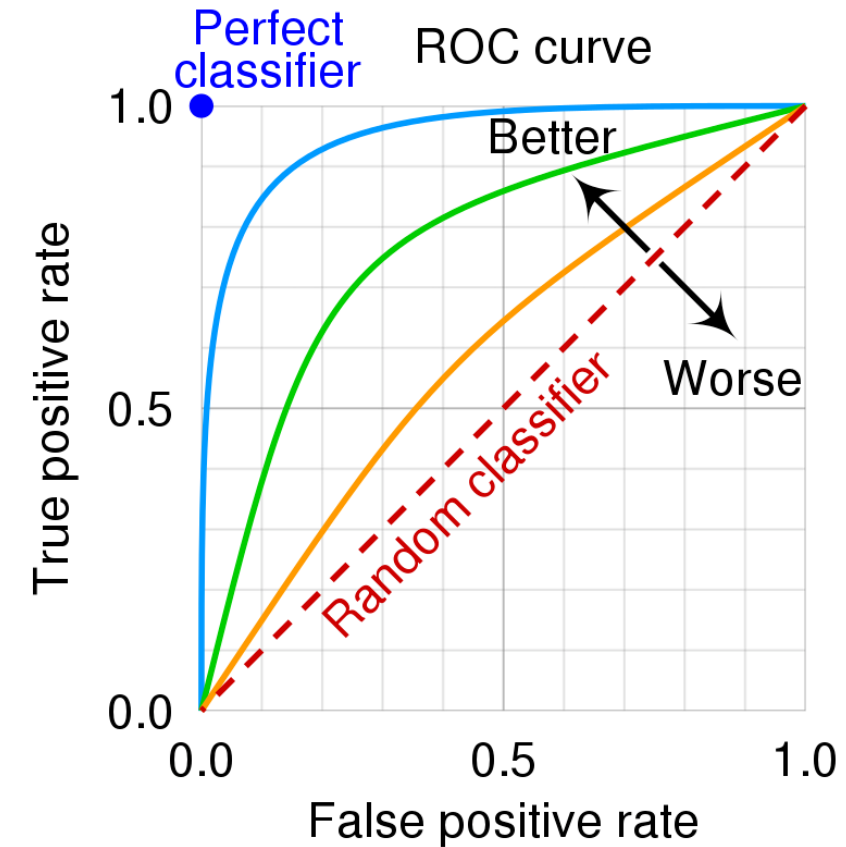
# Interpreting the quality of an AI/ML model



TEXAS A&M  
UNIVERSITY

Sources: [1][2][3][4][5][6][7][8][9] view · talk · edit

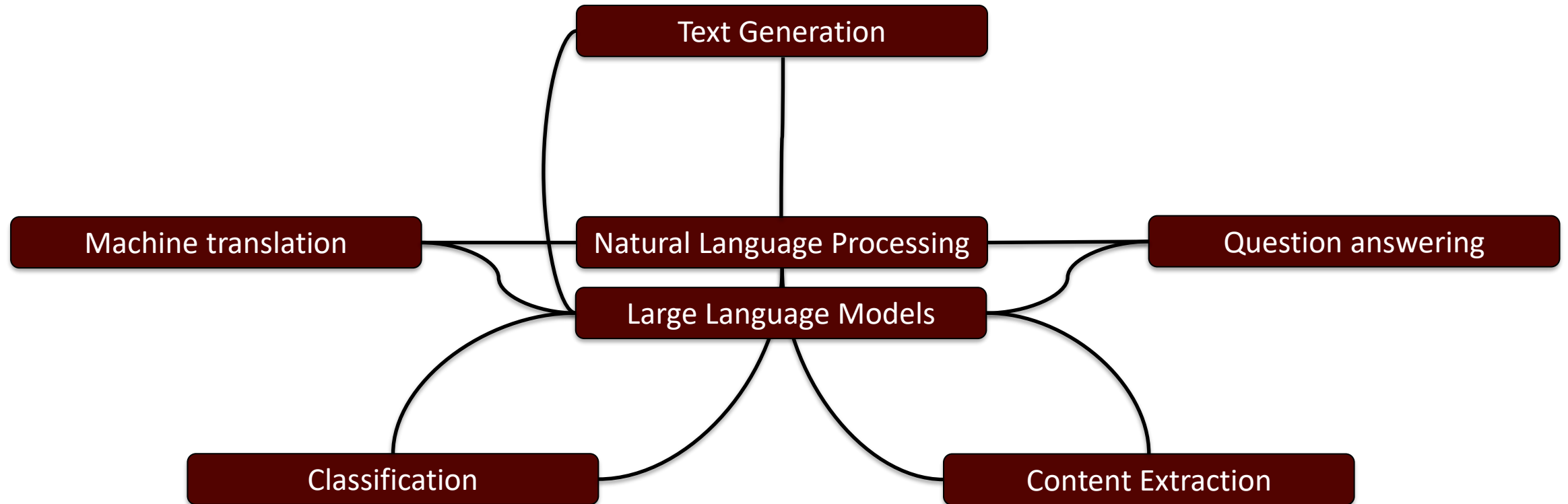
	Predicted condition			
	Positive (PP)	Negative (PN)		
Total population = P + N			Informedness, bookmaker informedness (BM) = TPR + TNR - 1	Prevalence threshold (PT) $= \frac{\sqrt{TPR \times FPR} - FPR}{TPR - FPR}$
Actual condition	Positive (P)	False negative (FN), type II error, miss, underestimation	True positive rate (TPR), recall, sensitivity (SEN), probability of detection, hit rate, power $= \frac{TP}{P} = 1 - FNR$	False negative rate (FNR), miss rate $= \frac{FN}{P} = 1 - TPR$
	Negative (N)	True negative (TN), correct rejection	False positive rate (FPR), probability of false alarm, fall-out $= \frac{FP}{N} = 1 - TNR$	True negative rate (TNR), specificity (SPC), selectivity $= \frac{TN}{N} = 1 - FPR$
Prevalence $= \frac{P}{P + N}$	Positive predictive value (PPV), precision $= \frac{TP}{PP} = 1 - FDR$	False omission rate (FOR) $= \frac{FN}{PN} = 1 - NPV$	Positive likelihood ratio (LR+) $= \frac{TPR}{FPR}$	Negative likelihood ratio (LR-) $= \frac{FNR}{TNR}$
Accuracy (ACC) $= \frac{TP + TN}{P + N}$	False discovery rate (FDR) $= \frac{FP}{PP} = 1 - PPV$	Negative predictive value (NPV) = $\frac{TN}{PN} = 1 - FOR$	Markedness (MK), deltaP ( $\Delta p$ ) = PPV + NPV - 1	Diagnostic odds ratio (DOR) $= \frac{LR+}{LR-}$
Balanced accuracy (BA) $= \frac{TPR + TNR}{2}$	F <sub>1</sub> score $= \frac{2PPV \times TPR}{PPV + TPR} = \frac{2TP}{2TP + FP + FN}$	Fowlkes–Mallows index (FM) $= \sqrt{PPV \times TPR}$	Matthews correlation coefficient (MCC) $= \frac{\sqrt{TPR \times TNR \times PPV \times NPV}}{\sqrt{FNR \times FPR \times FOR \times FDR}}$	Threat score (TS), critical success index (CSI), Jaccard index $= \frac{TP}{TP + FN + FP}$



# Natural language processing



TEXAS A&M  
UNIVERSITY



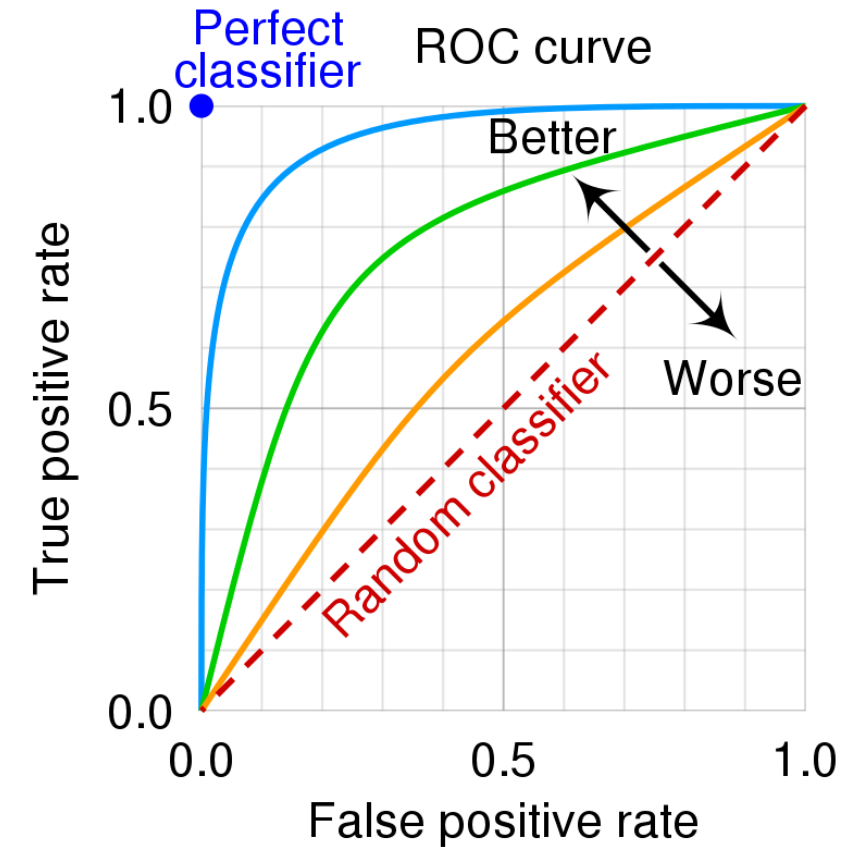
# Interpreting the quality of an AI/ML model



TEXAS A&M  
UNIVERSITY

Sources: [1][2][3][4][5][6][7][8][9] view · talk · edit

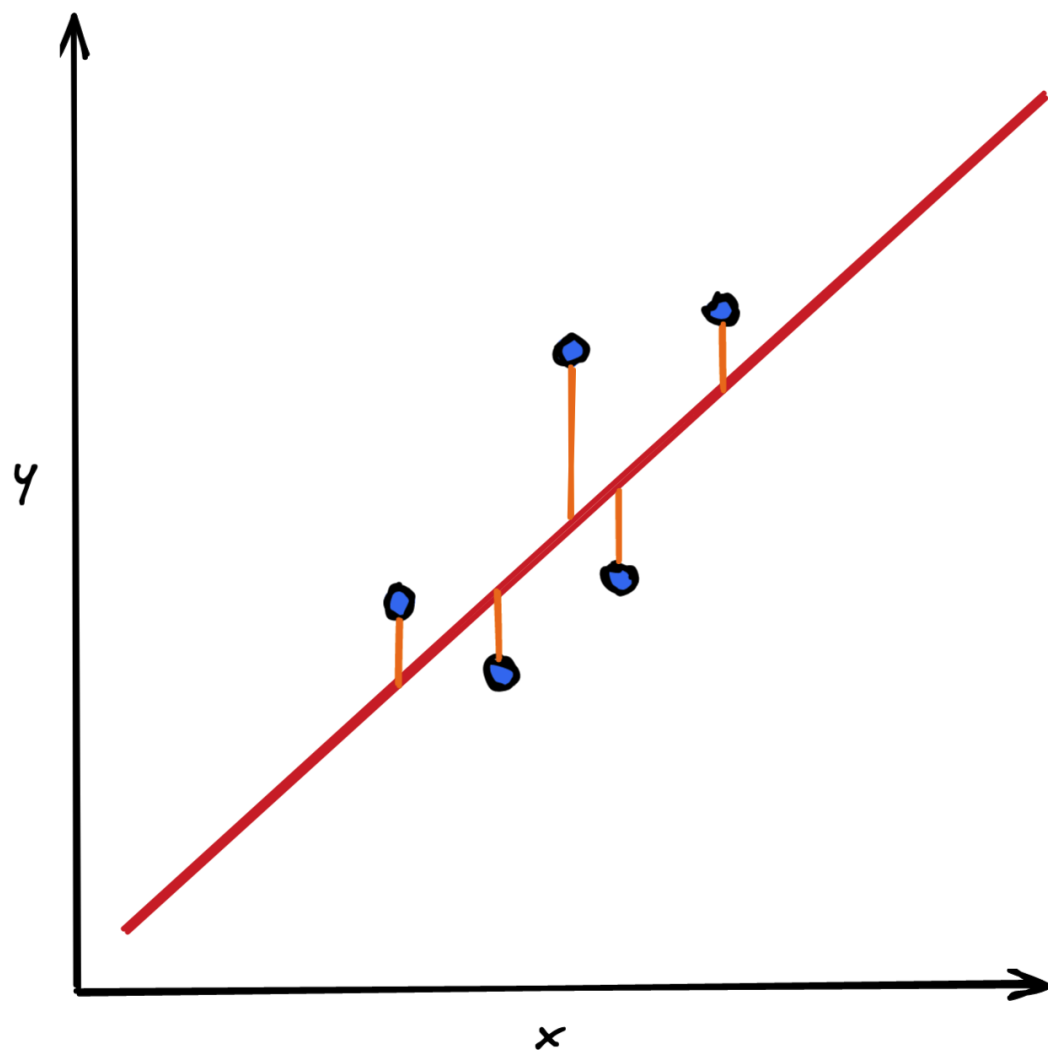
	Predicted condition			
	Positive (PP)	Negative (PN)		
Total population = P + N			Informedness, bookmaker informedness (BM) = TPR + TNR - 1	Prevalence threshold (PT) $= \frac{\sqrt{TPR \times FPR} - FPR}{TPR - FPR}$
Actual condition	Positive (P)	False negative (FN), type II error, miss, underestimation	True positive rate (TPR), recall, sensitivity (SEN), probability of detection, hit rate, power $= \frac{TP}{P} = 1 - FNR$	False negative rate (FNR), miss rate $= \frac{FN}{P} = 1 - TPR$
	Negative (N)	True negative (TN), correct rejection	False positive rate (FPR), probability of false alarm, fall-out $= \frac{FP}{N} = 1 - TNR$	True negative rate (TNR), specificity (SPC), selectivity $= \frac{TN}{N} = 1 - FPR$
Prevalence $= \frac{P}{P + N}$	Positive predictive value (PPV), precision $= \frac{TP}{PP} = 1 - FDR$	False omission rate (FOR) $= \frac{FN}{PN} = 1 - NPV$	Positive likelihood ratio (LR+) $= \frac{TPR}{FPR}$	Negative likelihood ratio (LR-) $= \frac{FNR}{TNR}$
Accuracy (ACC) $= \frac{TP + TN}{P + N}$	False discovery rate (FDR) $= \frac{FP}{PP} = 1 - PPV$	Negative predictive value (NPV) = $\frac{TN}{PN} = 1 - FOR$	Markedness (MK), deltaP ( $\Delta p$ ) = PPV + NPV - 1	Diagnostic odds ratio (DOR) $= \frac{LR+}{LR-}$
Balanced accuracy (BA) $= \frac{TPR + TNR}{2}$	F <sub>1</sub> score $= \frac{2PPV \times TPR}{PPV + TPR} = \frac{2TP}{2TP + FP + FN}$	Fowlkes–Mallows index (FM) $= \sqrt{PPV \times TPR}$	Matthews correlation coefficient (MCC) $= \frac{\sqrt{TPR \times TNR \times PPV \times NPV}}{\sqrt{FNR \times FPR \times FOR \times FDR}}$	Threat score (TS), critical success index (CSI), Jaccard index $= \frac{TP}{TP + FN + FP}$



# Evaluating a regressive model



TEXAS A&M  
UNIVERSITY



Metric	Benefits	Limitations
<b>R-squared (<math>R^2</math>)</b>	<ul style="list-style-type: none"><li>- High values = a good fit</li><li>- a measure of the proportion of variance</li></ul>	<ul style="list-style-type: none"><li>- Influenced by the sample size &amp; # of predictors</li><li>- May not be reliable when there are outliers or non-linear relationships in the data</li></ul>
<b>Adjusted R-squared (<math>R^2_{adj}</math>)</b>	<ul style="list-style-type: none"><li>- Similar to <math>R^2</math></li><li>Takes into account the number of predictors</li><li>- Provides a more accurate assessment of the model's performance when there are multiple predictors</li></ul>	<ul style="list-style-type: none"><li>- May be less informative when there are only a few predictors in the model</li></ul>
<b>Mean Squared Error (MSE)</b>	<ul style="list-style-type: none"><li>- Measures the average squared difference between the predicted and actual values</li></ul>	<ul style="list-style-type: none"><li>- Provides a measure of the magnitude of the errors in the model</li></ul>
<b>Mean Absolute Error (MAE)</b>	<ul style="list-style-type: none"><li>- Measures the average absolute difference between the predicted and actual values</li></ul>	<ul style="list-style-type: none"><li>- Provides a measure of the magnitude of the errors in the model</li></ul>
<b>Root Mean Squared Error (RMSE)</b>	<ul style="list-style-type: none"><li>- Measures the square root of the average squared difference between the predicted and actual values</li></ul>	<ul style="list-style-type: none"><li>- Provides a measure of the magnitude of the errors in the model</li></ul>