JAVIER DUARTE

US ATLAS ML TRAINING

JULY 25, 2023

Machine learning
in
particle physics

**Simulation-based inference**
- Unfolding
- Domain adaptation
- Parameter estimation
- BSM physics
- Diff-erentiable simulation

**Anomaly detection**
- Auto-encoders
- Normalizing flows

**Generative modeling / density estimation**
- Generative adversarial networks
- Gaussian processes
- Phase space generation
- Mixture models
- Diffusion models

**Decorrelation methods**
- Quantile regression
- Adversarial training

**Regression**
- Symbolic regression
- Function approx-imation
- Lattice guage theory
- Parton distribution functions
- Matrix elements
- Recasting
- Calibration
- Pileup

**Uncertainty quantification**
- Interpretability
- Estimation
- Mitigation
- Uncertainty-aware learning

**Representations/Architectures**
- Event images
- Jet images
- Sequences
- Trees
- Graphs
- Sets (point clouds)
- Equivariant models
- Physics-inspired

**Classification**
- Param-etrized classifiers
- Un-supervised
- Weak/semi-supervised

**Fast inference**
- Knowledge distillation
- Firmware/software
- Deployment
- Hardware-aware learning

**Learning strategies**
- Optimal transport
- Regular-ization
- Attention
- Feature ranking
- Quantum machine learning
- Reinforce-ment learning
- Hyper-parameter opti-mization

**Targets**
- Jet tagging
- BSM physics
- Particle identification
- Cosmology, astro-, and cosmic-ray physics
- Neutrino detectors
- Direct dark matter detectors

Full course:
[jduarte.physics.ucsd.edu/phys139_239](jduarte.physics.ucsd.edu/phys139_239)

[iml-wg.github.io/HEPML-LivingReview](iml-wg.github.io/HEPML-LivingReview)
[github.com/jmduarte/Nomological_Net_ML_Particle_Physics](github.com/jmduarte/Nomological_Net_ML_Particle_Physics)
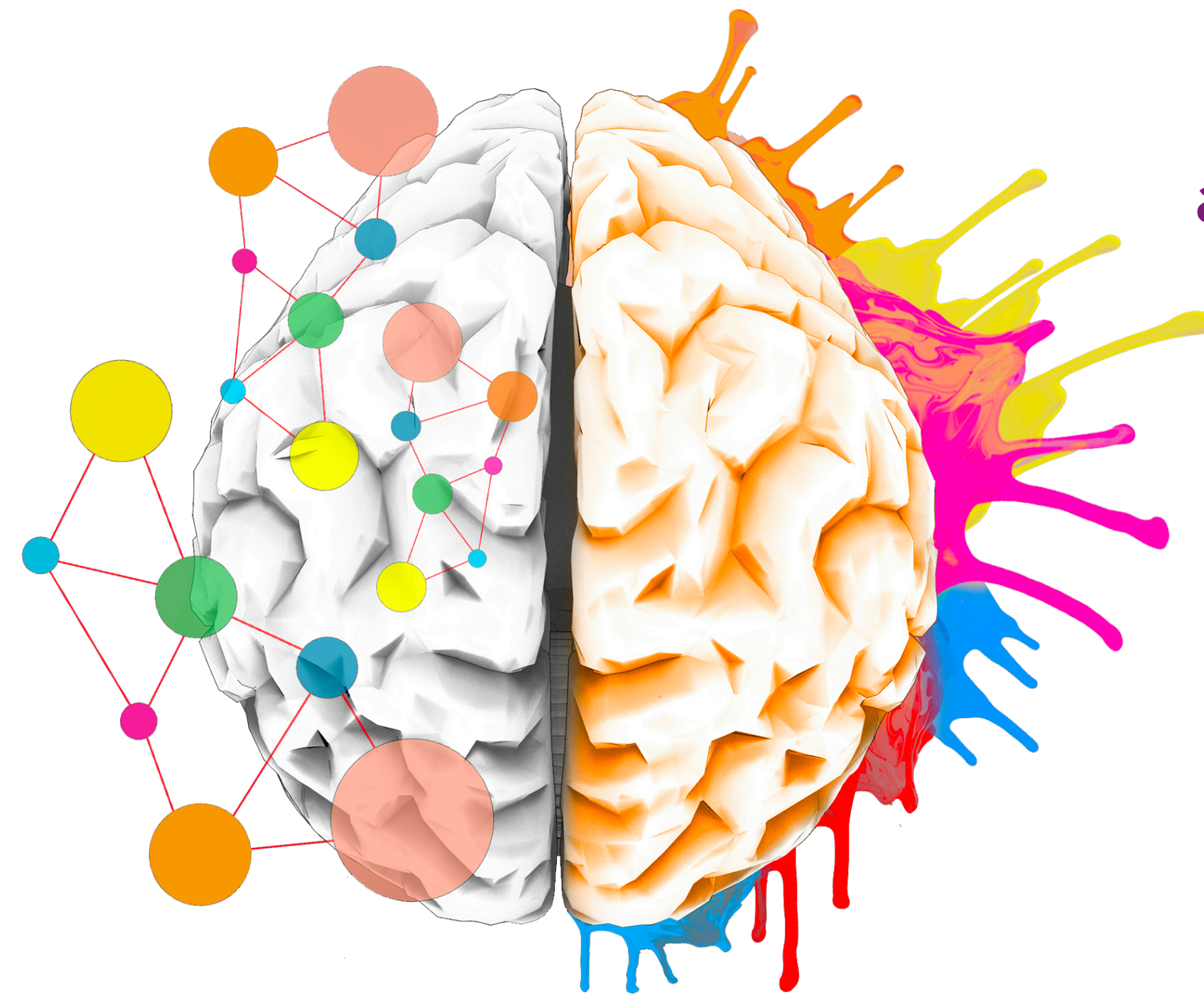
# I. BASICS
# II. DATA REPRESENTATIONS & SYMMETRIES
# III. ANOMALY DETECTION
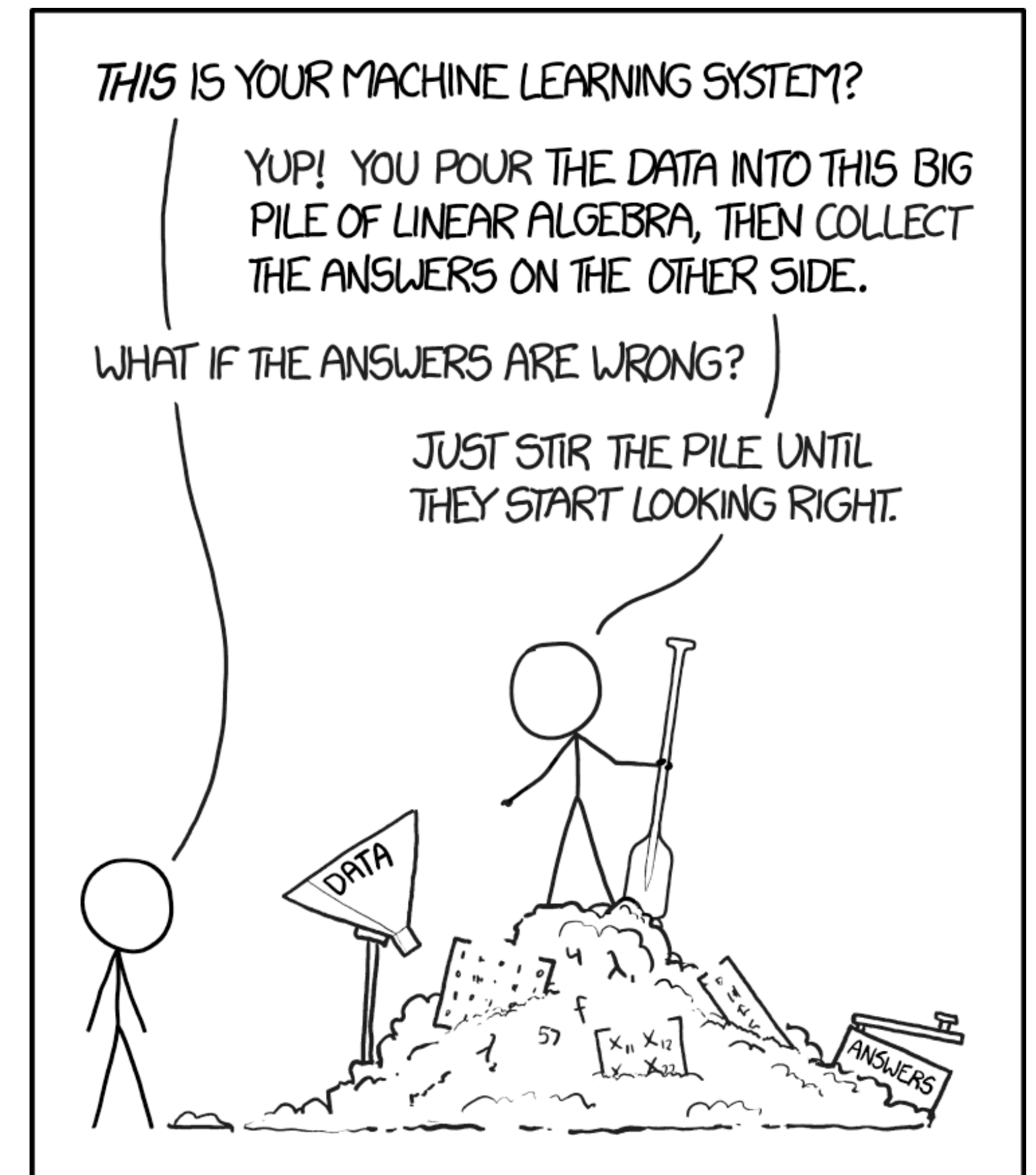# IV. GENERATIVE MODELING
# V. SUMMARY & OUTLOOK

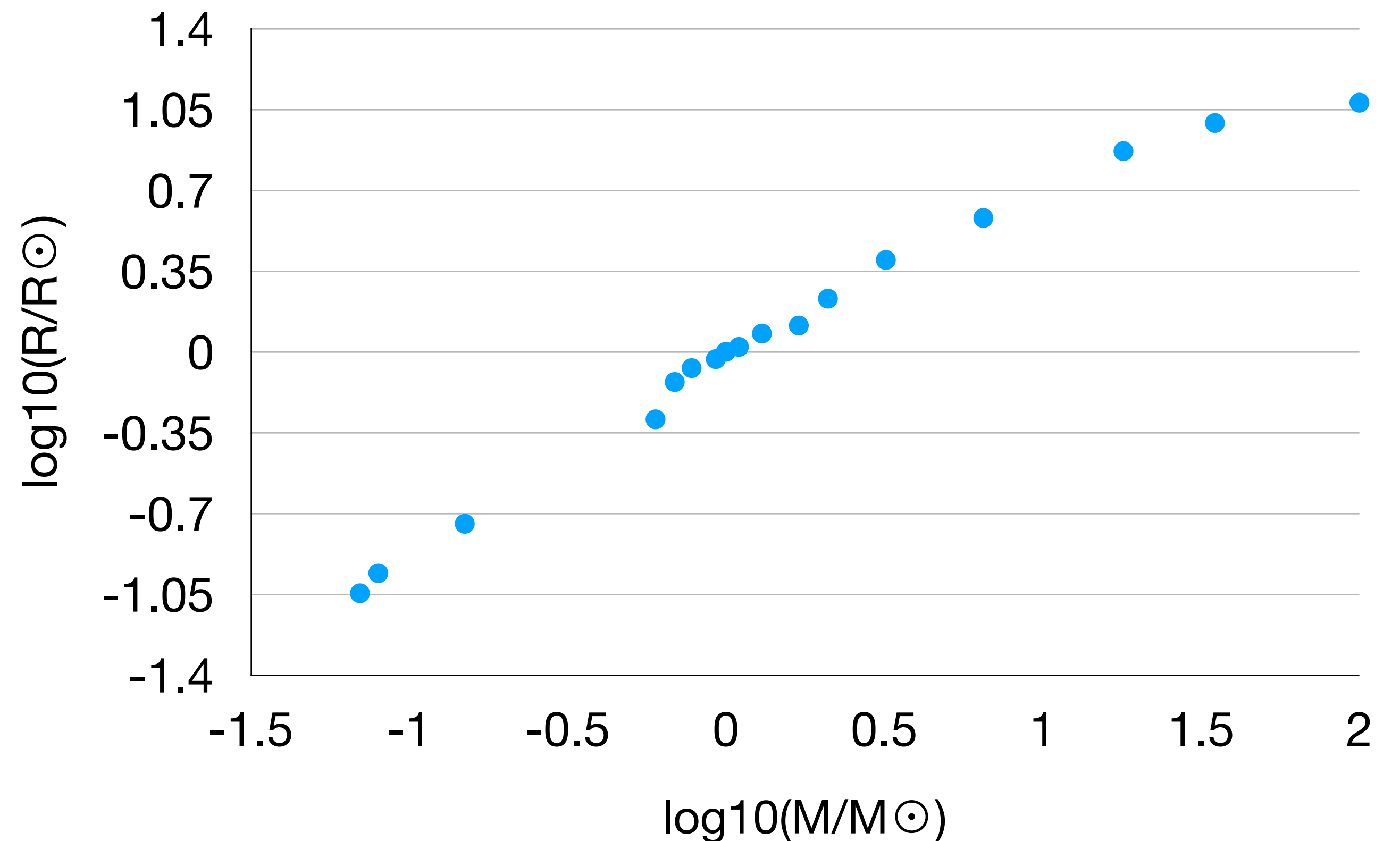▸ **Science** and **art** of learning automatically from data and experience

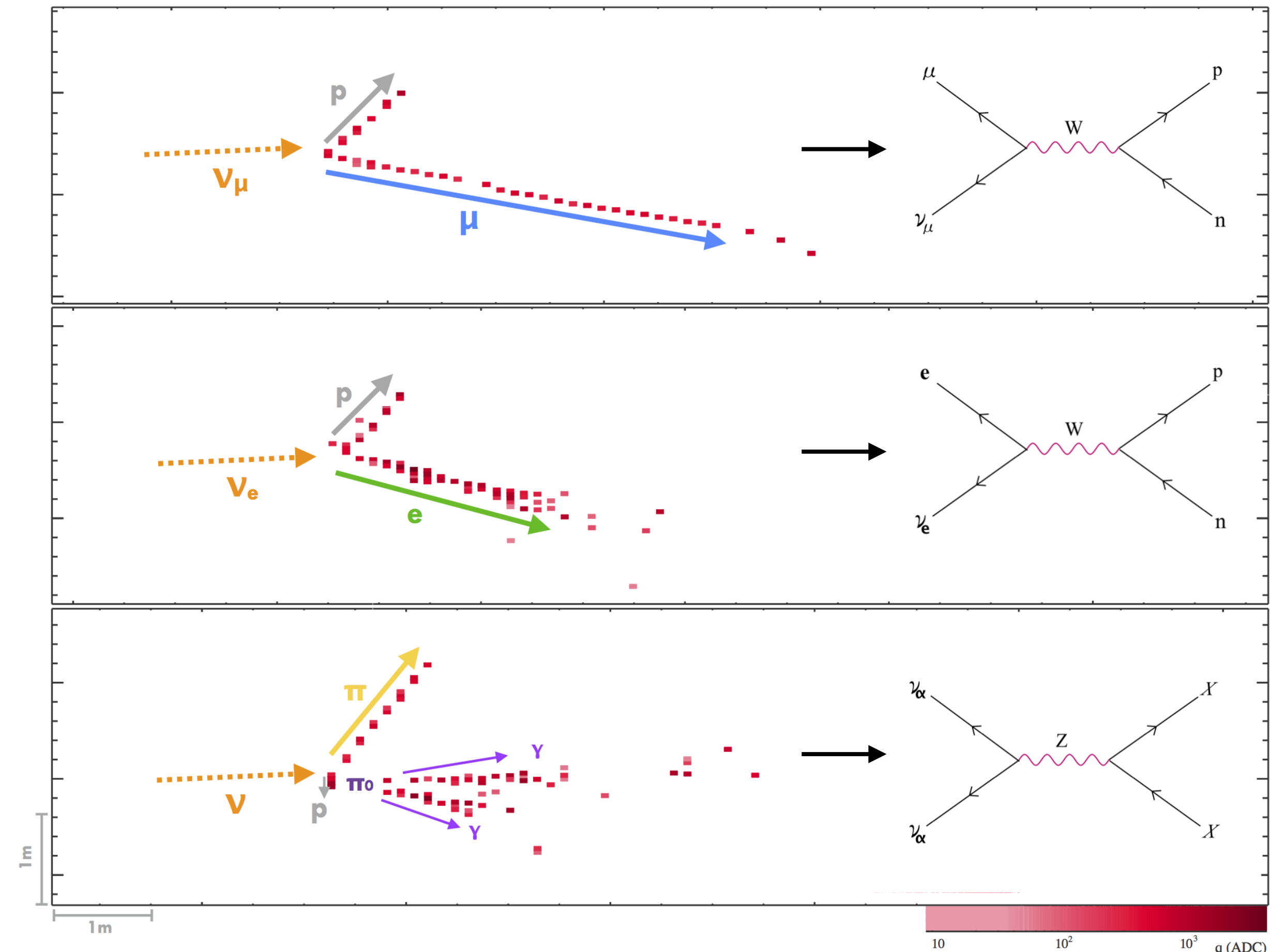Also, a lot of calculus, linear algebra, statistics, group theory, …



THIS IS YOUR MACHINE LEARNING SYSTEM?

YUP! YOU POUR THE DATA INTO THIS BIG PILE OF LINEAR ALGEBRA, THEN COLLECT THE ANSWERS ON THE OTHER SIDE.

WHAT IF THE ANSWERS ARE WRONG?

JUST STIR THE PILE UNTIL THEY START LOOKING RIGHT.

▸ Large overlap with data mining:
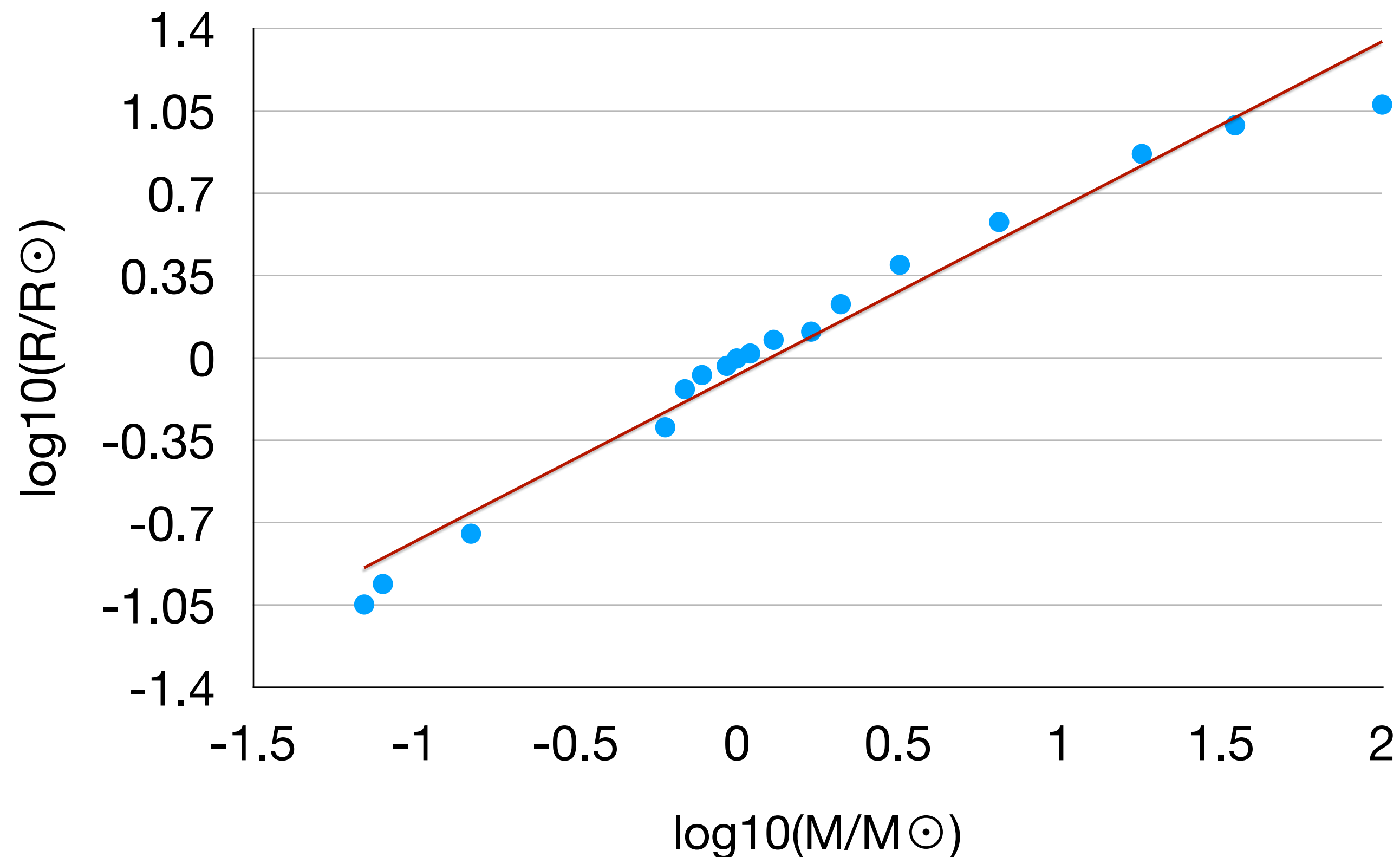
  ▸ ML focuses on algorithms, DM on discovering patterns

▸ Learn a function $f : X \rightarrow Y$ from an input space $X$ (observations) to an output space Y (targets), using a set of labeled examples $(x_1, y_1), (x_2, y_2), \ldots, (x_N, y_N)$.

▸ Example 1: Predict stellar radius given stellar mass

▸ Learn a function $f : X \to Y$ from an input space $X$ (observations) to an output space Y (targets), using a set of labeled examples $(x_1, y_1), (x_2, y_2), \ldots, (x_N, y_N)$.

▸ Example 2: Classify images of neutrino interactions



arXiv:1604.01444

▸ Learn a function $f : X \to Y$ from an input space $X$ (observations) to an output space Y (targets), using a set of labeled examples $(x_1, y_1), (x_2, y_2), \ldots, (x_N, y_N)$.

▸ Example 3: Reduce noise in a time-series trace to identify a gravitational wave signal

arXiv:1711.03121

▸ Learn a function $f : X \to Y$ from an input space $X$ (observations) to an output space Y (targets), using a set of labeled examples$(x_1, y_1), (x_2, y_2), \ldots, (x_N, y_N)$.

▸ Example 4: Estimate particle momentum, charge, type, etc. from detector hits



$$\to \begin{pmatrix} E \\ p_x \\ p_y \\ p_z \end{pmatrix}, q, \text{type}, p_{\text{pileup}}, \cdots$$

arXiv:2101.08578

‣ Collect a **labeled** training set (supervision)

    ‣ Often requires **simulation** where the "ground truth" is known



‣ Train a model using a learning algorithm (find patterns in the data)

$f(x_i | w)$

Error

$y_i$

log10(R/R⊙)

log10(M/M⊙)

$x_i$

▸ Linear model:

$$f(x | w) = w^\mathsf{T} x \quad (w \in \mathbb{R}^{D+1})$$

▸ How do we select the parameters w?

▸ We want $y_i \approx f(x_i | w)$

▸ Squared loss: $L(y, y') = (y - y')^2$

(Least squares)

Learning objective: $\arg\min_w \sum_{i=1}^{N} L(y_i, f(x_i | w)) = \arg\min_w \sum_{i=1}^{N} (y_i - w^\mathsf{T} x_i)^2$

▸ In supervised learning, we want to optimize the objective

$$l(w) = \sum_{i=1}^{N} L(y_i, f(x_i \mid w))$$

▸ For linear regression, there is a closed-form solution, but in general?

▸ We need an optimization algorithm to find the optimal (or just "good") w

▸ Set $w(t = 0)$ to some values (e.g., $w(0) = 0$ or some random value)

▸ At iteration $t$,

    ▸ Compute the gradient $\nabla_w l(w(t))$: direction of steepest increase of $l(w)$ at $w(t)$

    ▸ Take a small step in the opposite direction:

$$w(t + 1) = w(t) - \eta \nabla_w l(w(t))$$

Step size / learning rate

$l(w)$

$w^{(0)}$

$w^{(1)}$

▸ Fitting the training dataset perfectly (error = 0) does not necessarily mean the model will work well on new test data!



Linear fit: ok on both training and testing

Polynomial fit (degree 4): excellent on training, bad on testing

▸ If $L$ is the squared loss, we can decompose the expected test error:

$$\mathbb{E}\left[L_P(f(x\,|\,w_S)\right] = \mathbb{E}_S\mathbb{E}_{(x,y)\sim P(x,y)}\left[L(y, f(x\,|\,w_S))\right]$$

$$= \mathbb{E}_{(x,y)\sim P(x,y)}\left[\underbrace{\mathbb{E}_S\left[(f(x\,|\,w_S) - F(x))^2\right]}_{\text{Variance}} + \underbrace{(F(x) - y)^2}_{\text{(Squared) bias}}\right]$$

▸ where $F(x) = \mathbb{E}_S\left[f(x\,|\,w_S)\right]$ is the average prediction of our model over different possible training datasets

▸ **Variance**: difference in predictions when training on different datasets

▸ **Bias**: difference from ground truth

# OVERFITTING VS. UNDERFITTING

▸ Overfitting implies high variance (unstable model class)

  ▸ Variance increases with model complexity

  ▸ Variance decreases with more training data

▸ Underfitting implies high bias

  ▸ Even with no variance, model class has high error

  ▸ Underfitting happens whenever model complexity is too low



Degree 5



Degree 0

features / embedding of $x$

▸ Replace our input vector $x$ with some $\phi(x)$ to make our model more expressive

▸ For example, if $\phi(x) = (1, x, x^2)$ then our model becomes:

$$f(x \mid w) = w^\mathsf{T}\phi(x) = w_0 + w_1 x + w_2 x^2$$



- The model is still **linear** in the parameters $w$!

- More expressive than a line $w_0 + w_1 x$, so the fit is better (i.e., training error is lower)

▸ Linear models on top of good features can yield excellent results

▸ More complex model classes (e.g., neural networks) have linear models as their basic building block

    ▸ NNs are "automatic featurizers"

Linear model: $f(x \mid w) = w^\mathsf{T} x$

Neural network: linear model after inputs are mapped to features through a nonlinear transformation

$$f(x \mid w_1, w_2) = w_2^\mathsf{T} \sigma(w_1^\mathsf{T} x)$$

▸ We only have a finite training dataset

▸ We cannot measure the true test error

▸ Simple model classes underfit ⎤

                                         Bias-variance tradeoff

▸ Complex model classes overfit ⎦

(but not so straightforward for deep neural networks!)

▸ **Goal**: Select the model class with the lowest test error

Original dataset

▸ Split the original dataset into a training and validation set

▸ Train model on the training set

▸ Evaluate on the validation set to estimate the test error

▸ Select the model class that gives the lowest estimated error

▸ Optionally, re-train the selected model class on the whole dataset (training + validation)

▸ **Issue**: we would like both training and validation sets to be as large as possible (so that the estimate is better), but they must not overlap!

▸ Split the original dataset into $k$ equal parts (e.g, $k = 5$)

▸ Train on the $k - 1$ parts and validate on the remaining one



Original dataset

▸ Repeat for every choice of the $k - 1$ parts and average the validation errors



▸ **Advantage**: use all data as validation to improve the estimate of the test error, at the cost of more computation ($k$ trainings)

▸ Training dataset: $S = \{(x_1, y_1), \ldots, (x_N, y_N)\}$ where $x \in \mathbb{R}^D$ and $y \in \mathbb{R}$

or $\phi(x)$ instead of $x$

▸ Model / hypothesis class: $f(x \,|\, w) = w^\mathsf{T} x$ (linear models)

▸ Loss function: $L(y, y') = (y - y')^2$ (squared loss)

▸ Optimization algorithm to minimize the learning objective:

$$\arg\min_{w} \sum_{i=1}^{N} L(y_i, f(x_i \,|\, w))$$

▸ Cross validation and model selection:

▸ Testing and deployment

**Important**: if a testing set is available, never use it to make decisions on the model!

UCSD

- High-level (expert) variables

- Ordered list of particles

- Images

- Set of particles

- Graph of particles

- Lorentz scalars/vectors

- Shallow neural network, boosted decision tree, …

- 1D convolutional neural network, recurrent neural network

- 2D convolutional neural network

- Deep set (energy flow network)

- Graph neural network

- Lorentz-equivariant network

TRANSFORMERS

▸ After "particle-flow reconstruction,"  can think of event as a collection of points in momentum space

▸ For jets (localized clusters of particles), dimensionality ($N_{particles} \sim 100,\ 4 + M$)

▸ Variable jet length requires:

  ▸ Preprocessing into another rep. (tab. data, jet images, …)

  ▸ Truncation to fixed size

  ▸ Graph NN

u,d or s jet

c or b jet

gluon jet

pileup jet

W or Z jet

Higgs jet

top jet

?

▸ Tabular data: use physics knowledge to preprocess jet information into a set of high-level features

▸ Substructure variable:

$$_1e_3^\beta = \sum_{1 \le i < j < k \le n_J} z_i z_j z_k \min\{\Delta R_{ij}^\beta, \Delta R_{ik}^\beta, \Delta R_{jk}^\beta\}$$

$$_2e_3^\beta = \sum_{1 \le i < j < k \le n_J} z_i z_j z_k \min\{\Delta R_{ij}^\beta \Delta R_{ik}^\beta, \Delta R_{ij}^\beta \Delta R_{jk}^\beta, \Delta R_{ik}^\beta \Delta R_{jk}^\beta\}$$

   ▸ jet mass

   ▸ energy correlation functions, e.g. $N_2^{\beta=1} = {}_2e_3^{\beta=1}/({}_1e_3^{\beta=1})^2$

MiniBooNE detector

Signal region

Veto region

Root node

S/B
52/48

< 100    PMT Hits?    ≥ 100

B
4/37

Branch node
(further branching)

S/B
48/11

< 0.2 GeV    Energy?    ≥ 0.2 GeV

Leaf nodes
(no further branching)

S/B
9/10

S
39/1

< 500 cm    Radius?    ≥ 500 cm

S
7/1

B
2/9

MiniBooNE: 1520 photomultiplier signals
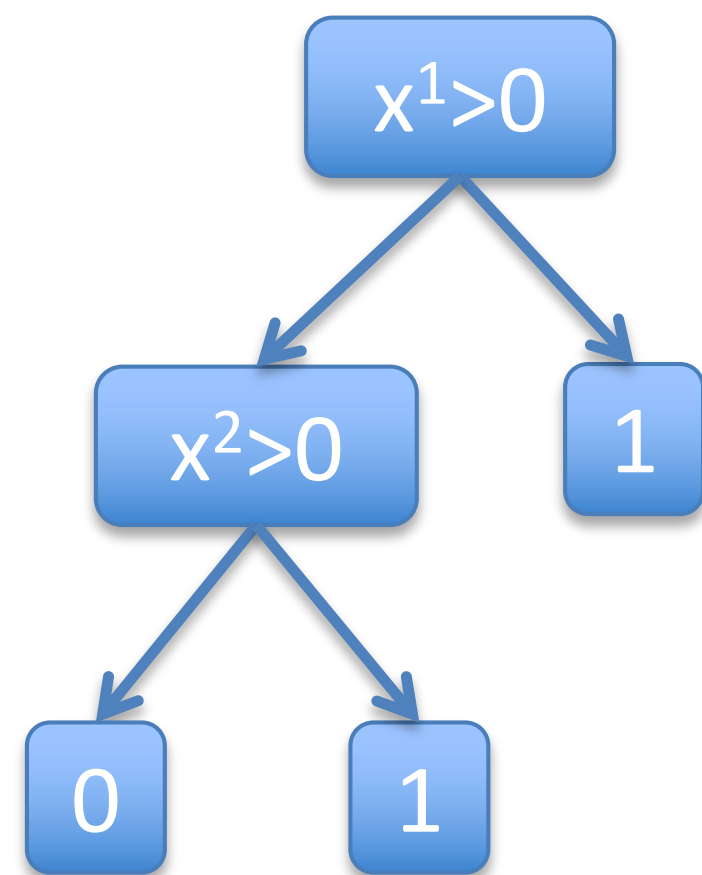Goal: separate $\nu_e$ and $\nu_\mu$ events
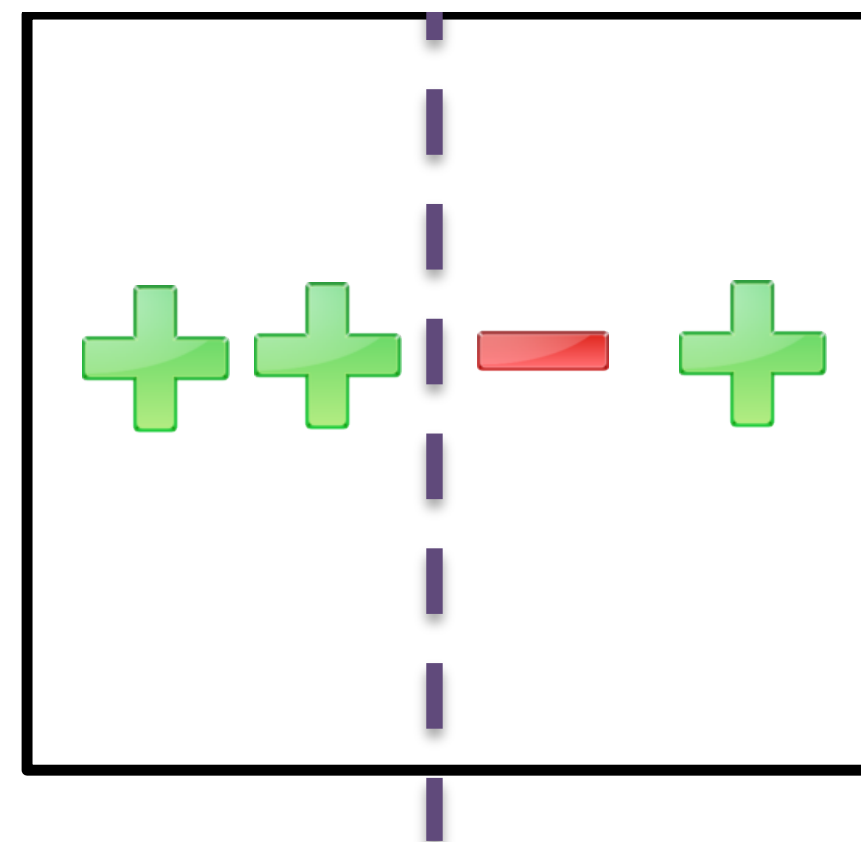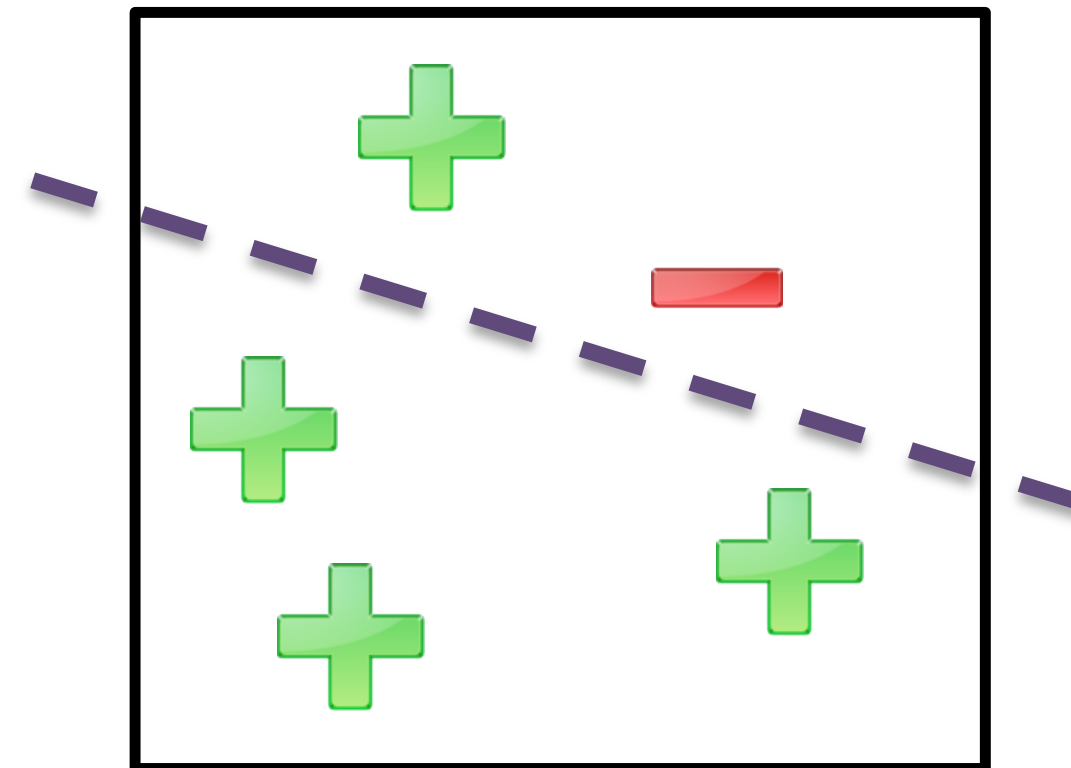
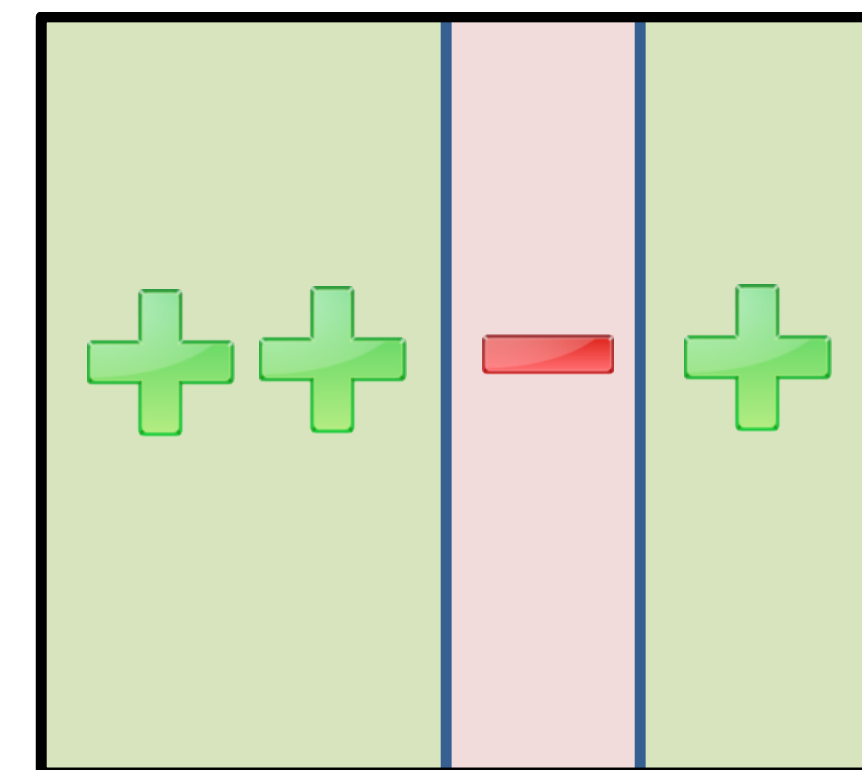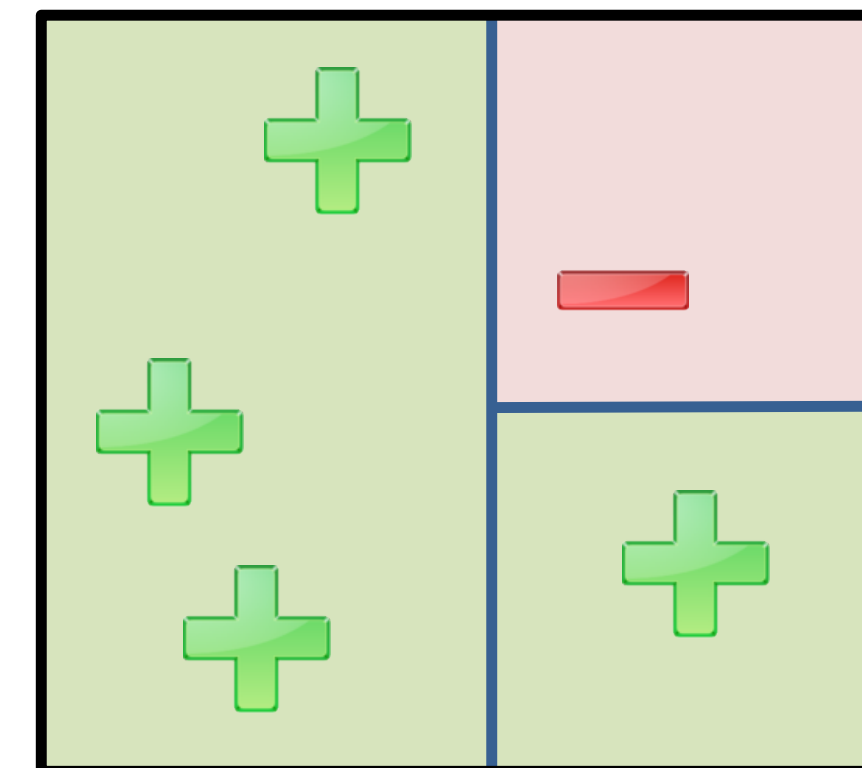$\nu_e$ CCQE

$\nu_e n \rightarrow p e^-$

$\nu_\mu$ CCQE

$\nu_\mu n \rightarrow p \mu^-$

▸ Leaf nodes classify events as either
signal ($\nu_e$) or background ($\nu_\mu$)

▸ Decision trees are nonlinear models!

▸ Examples:

No linear model
can achieve 0 error

Simple decision tree
can achieve 0 error

$x^1>0$

$x^2>0$

1

0

1

$x^1>0$

$x^1>1$
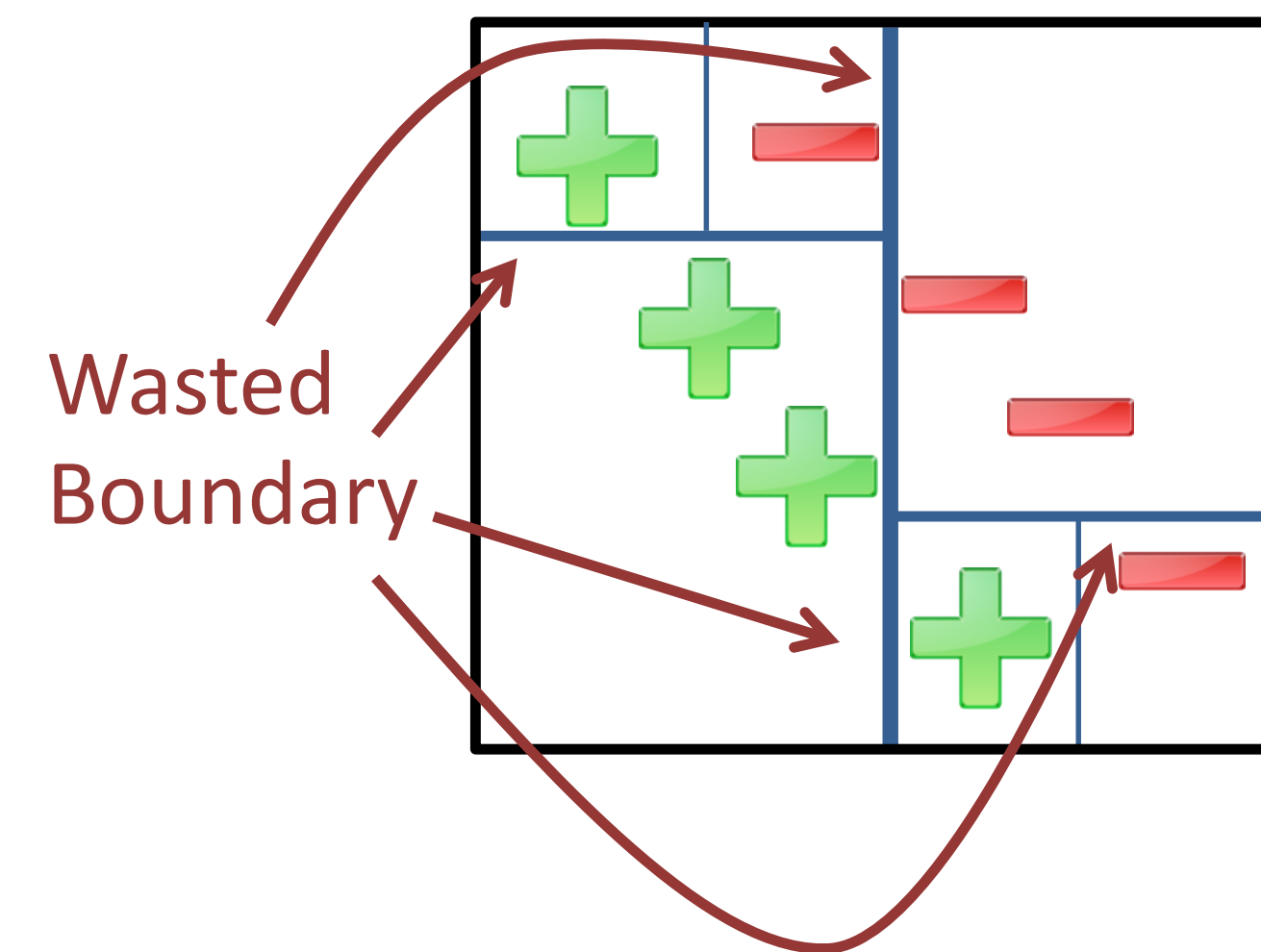
1

1

0

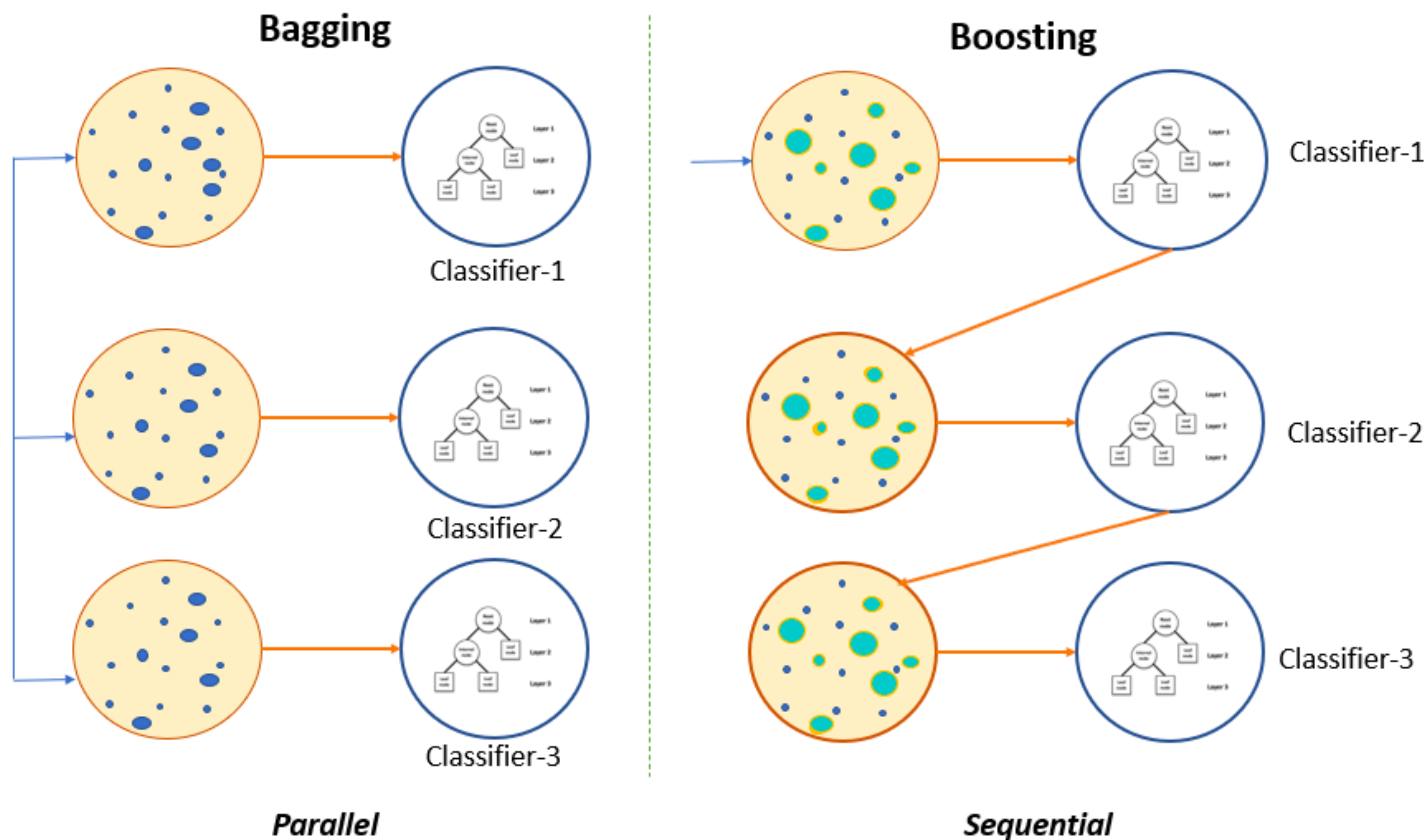▶ Decision trees are axis-aligned!

▶ Example:

Simple linear SVM can
easily find max margin
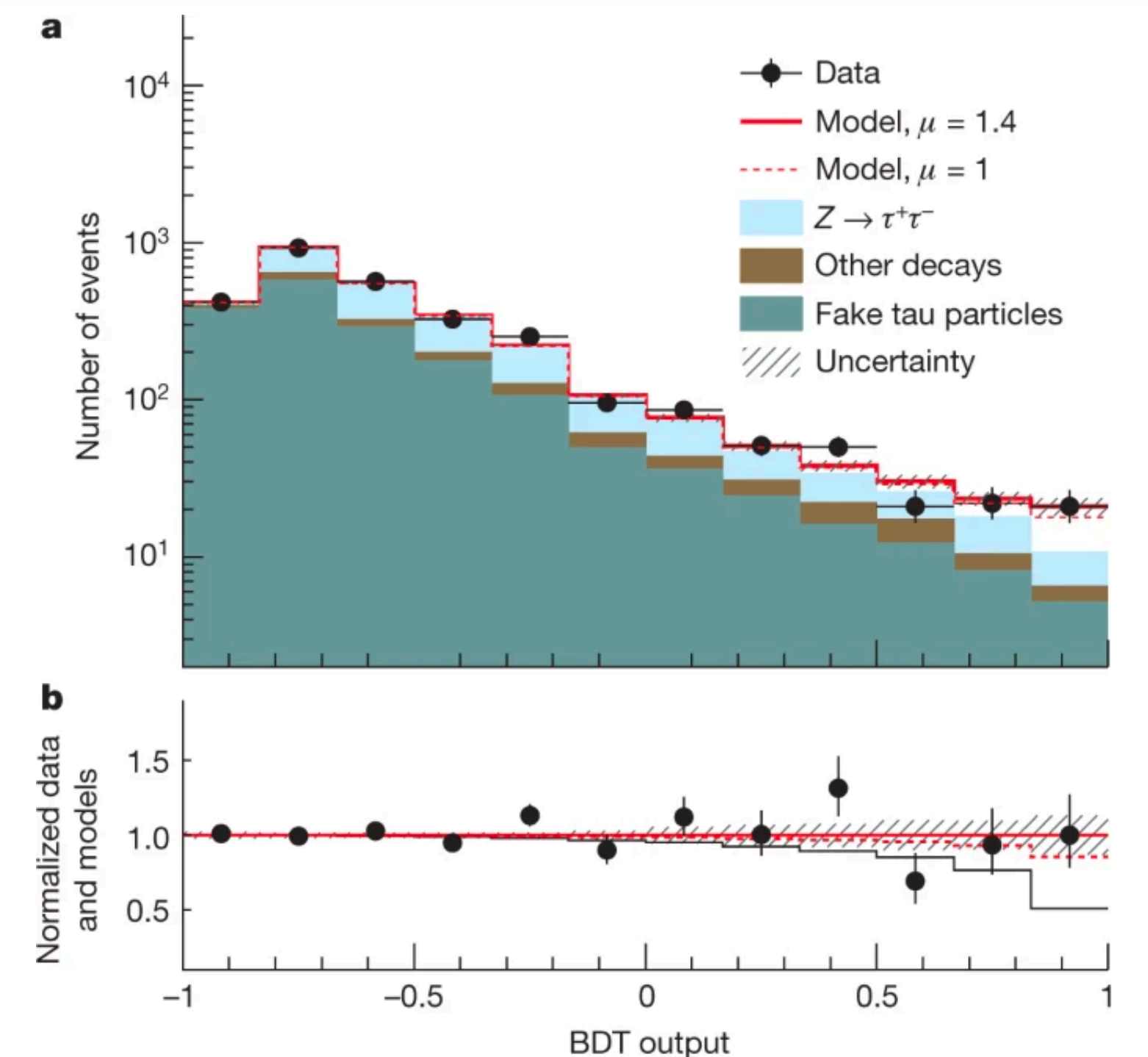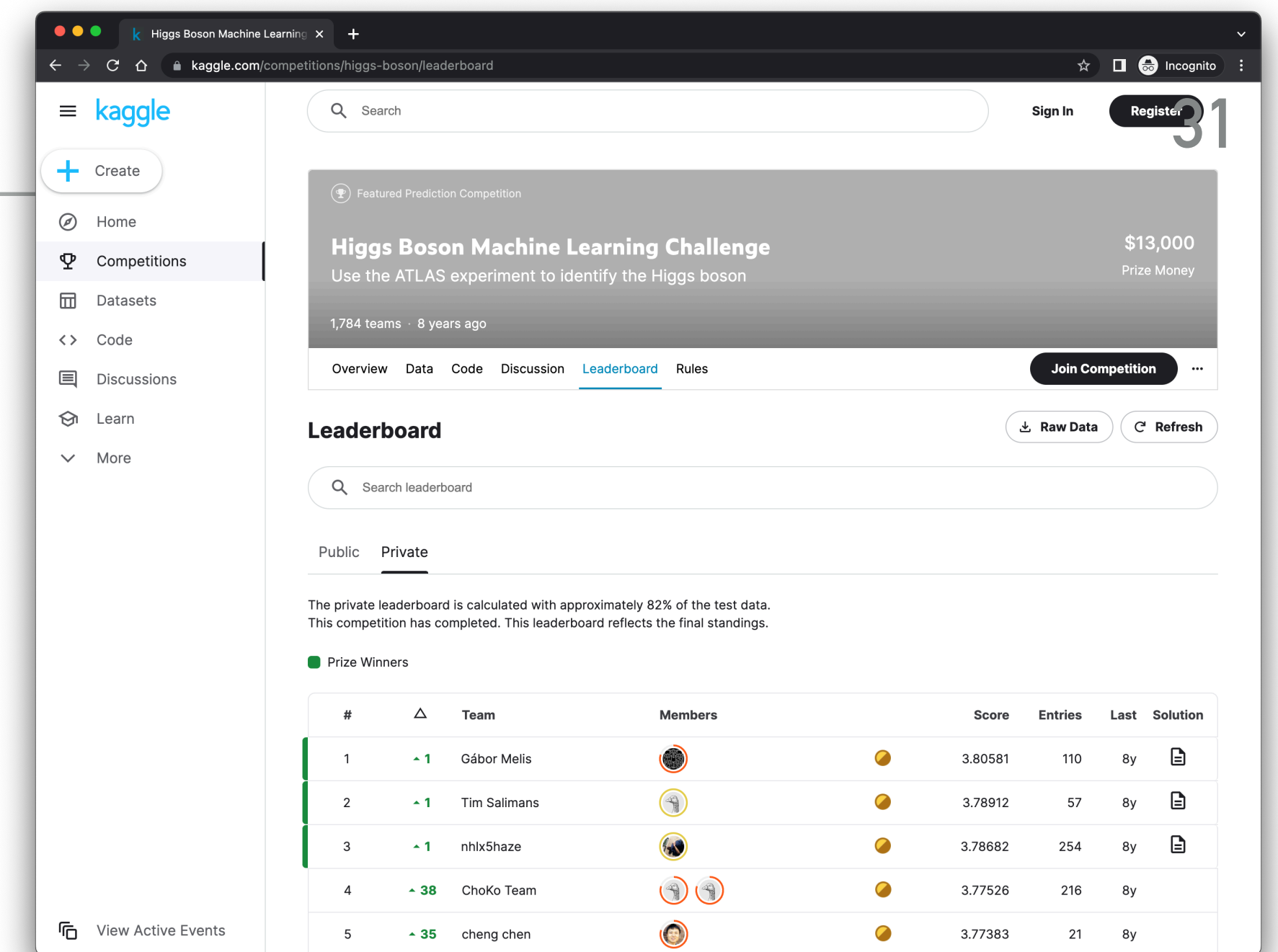
Decision trees require
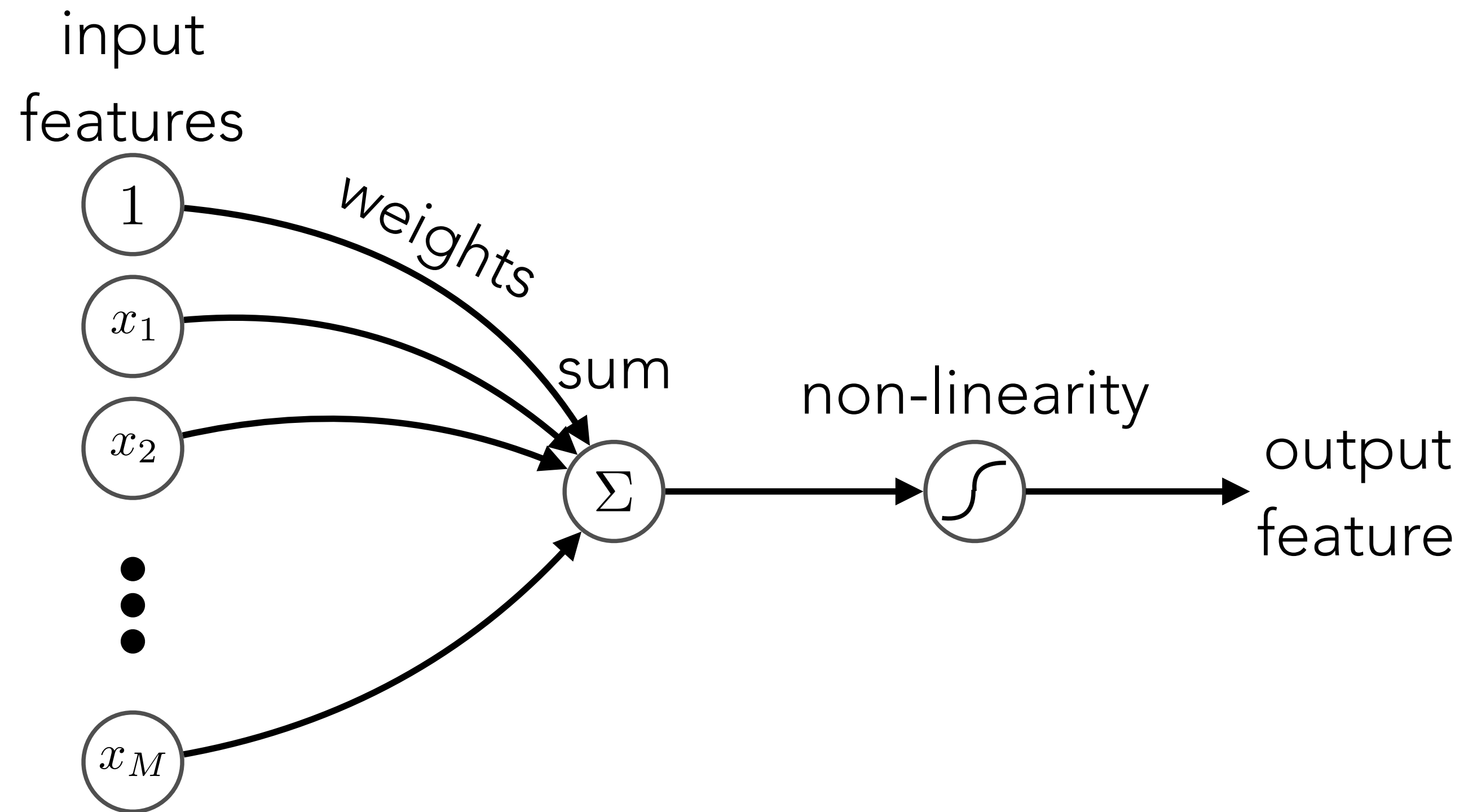complex axis-aligned
partitioning

Wasted
Boundary

▸ Bagging: reduce variance of weak learners

▸ Boosting: reduce bias of weak learners

# BOOSTED DECISION TREES IN THE WILD

▸ 1st place in Kaggle Higgs Boson Machine Learning Challenge [kaggle.com/competitions/higgs-boson]

  ▸ And many other uses at LHC, e.g. in Higgs boson discovery [10.1038/s41586-018-0361-2]

▸ Predicting critical temperature of a superconductor [10.1016/j.commatsci.2018.07.052]

▸ MiniBooNE neutrino event classification [10.1016/j.nima.2004.12.018]

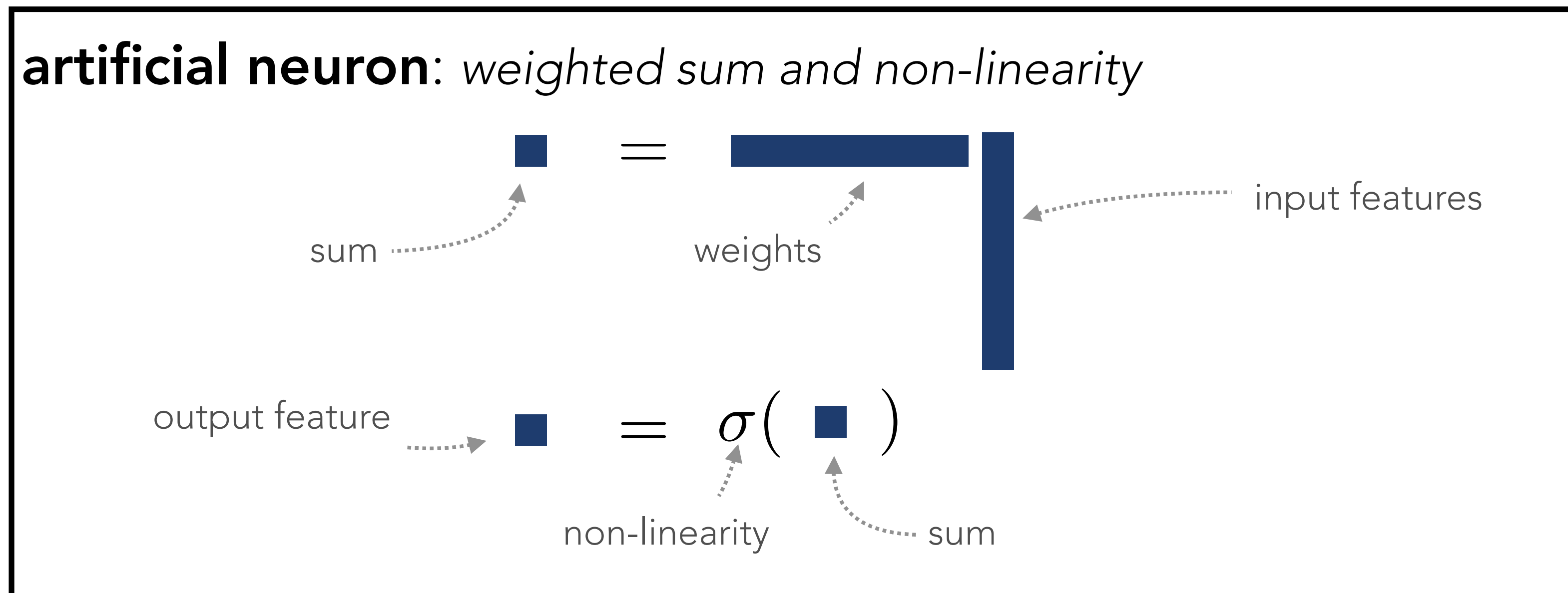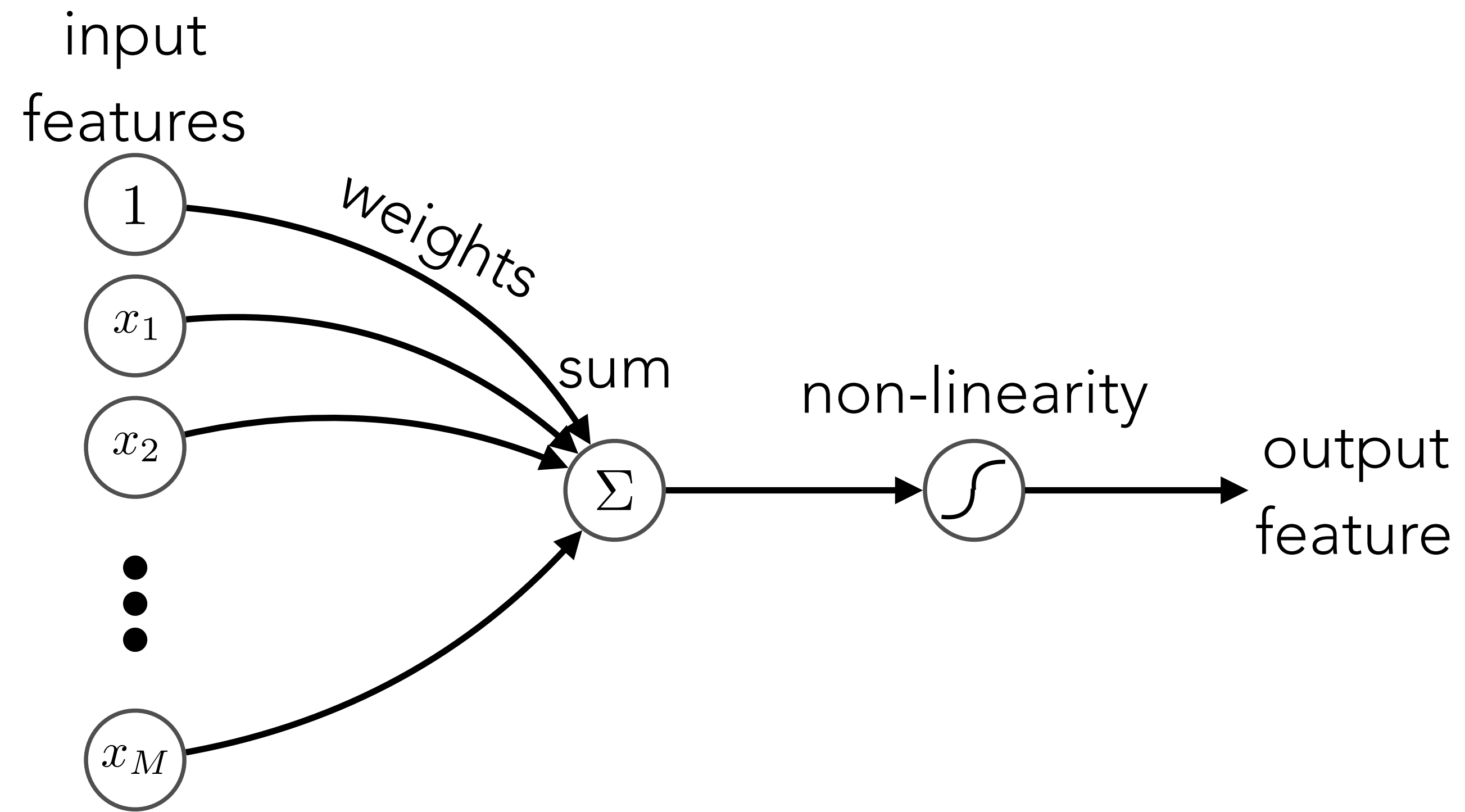▸ Observation of single top quark production at D0 [10.1103/PhysRevLett.103.092001]
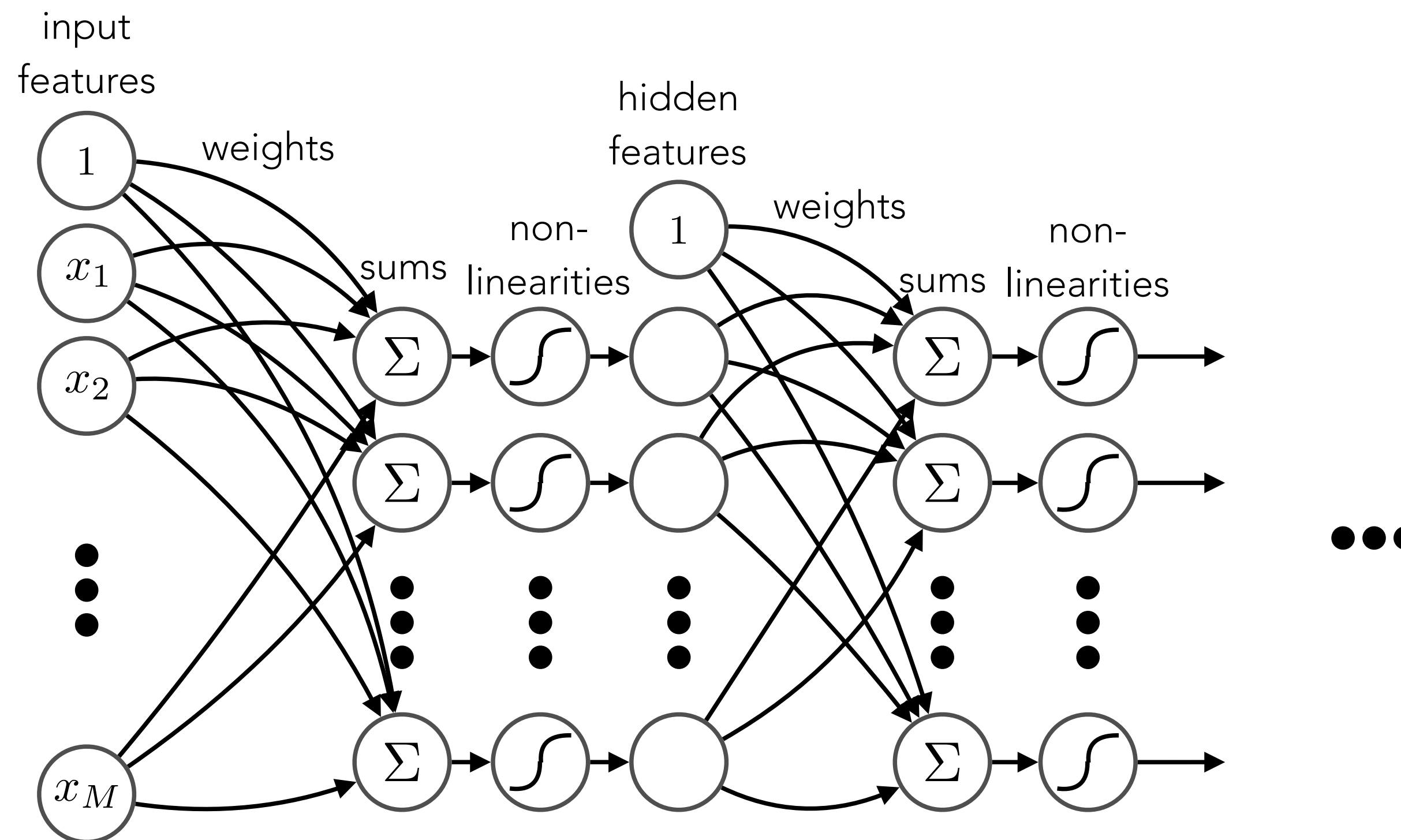
input
features



weights

sum

non-linearity

output
feature

1

$x_1$

$x_2$

$x_M$

$\Sigma$

**artificial neuron**: *weighted sum and non-linearity*

bias

input features

$$s = w_0 + w_1 x_1 + w_2 x_2 + \cdots + w_M x_M = \mathbf{w}^\mathsf{T}\mathbf{x}$$

sum

weights

$$h = \sigma(s)$$

output feature

non-linearity

sum

input
features

1

$x_1$

$x_2$

$x_M$

weights

sum

$\Sigma$

non-linearity

$\int$

output
feature

**artificial neuron**: *weighted sum and non-linearity*

■ = ■

sum          weights          input features

output feature   ■ = $\sigma($ ■ $)$

non-linearity          sum

input
features

weights

hidden
features

weights

non-
linearities

non-
linearities

sums

sums

$1$

$x_1$

$x_2$

$x_M$

$1$

$\Sigma$ $\int$

$\Sigma$ $\int$

$\Sigma$ $\int$

$\Sigma$ $\int$

$\Sigma$ $\int$

$\Sigma$ $\int$

**network**: *sequence of parallelized weighted sums and non-linearities*

$$\blacksquare = \sigma( \;\cdots\; \sigma( \;\blacksquare\; \sigma( \;\blacksquare\blacksquare\; ) \,)\;\cdots) $$

output

2nd weights

1st weights

input

Universal approximation theorem (informal). Given a function $y = f(x)$ and an $\epsilon > 0$, there exists a deep network $y = f_w(x)$ (of arbitrary width or depth) such that:

$$\sup_{x \in X} \|f(x) - f_w(x)\| < \epsilon$$



Note: This means that a network can *represent* any function, not that it can learn it! The "amount" of function a given network can represent is often called its expressive power.

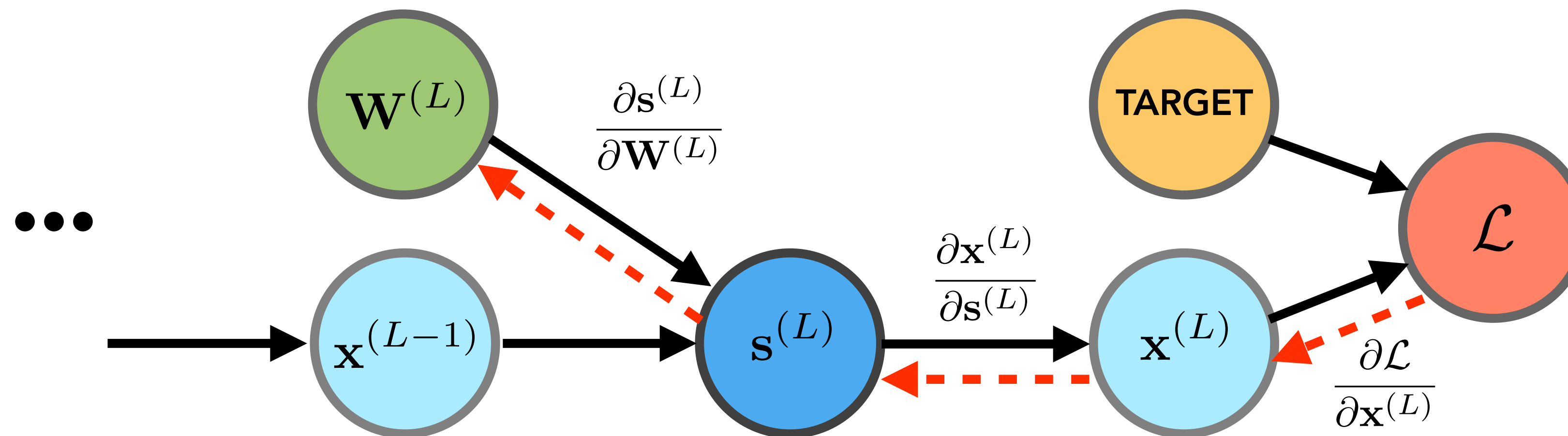We train deep networks using Maximum Likelihood Estimation (MLE): The last layer of a DNN is a softmax that outputs probabilities over classes:



x = input data

$$p_w(y \mid x) = \begin{bmatrix} 0.9 \\ 0.1 \\ \vdots \end{bmatrix}$$

w = vector containing all weights

y = label

We train the weights $w$ to maximize the log-likelihood of the data under our model:

$$L(w) = -\frac{1}{N} \sum_{i=1}^{N} \log p_w(y_i \mid x_i)$$

Negative log-likelihood loss (cross-entropy loss)

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}^{(L)}} = \frac{\partial \mathcal{L}}{\partial \mathbf{x}^{(L)}} \frac{\partial \mathbf{x}^{(L)}}{\partial \mathbf{s}^{(L)}} \frac{\partial \mathbf{s}^{(L)}}{\partial \mathbf{W}^{(L)}}$$
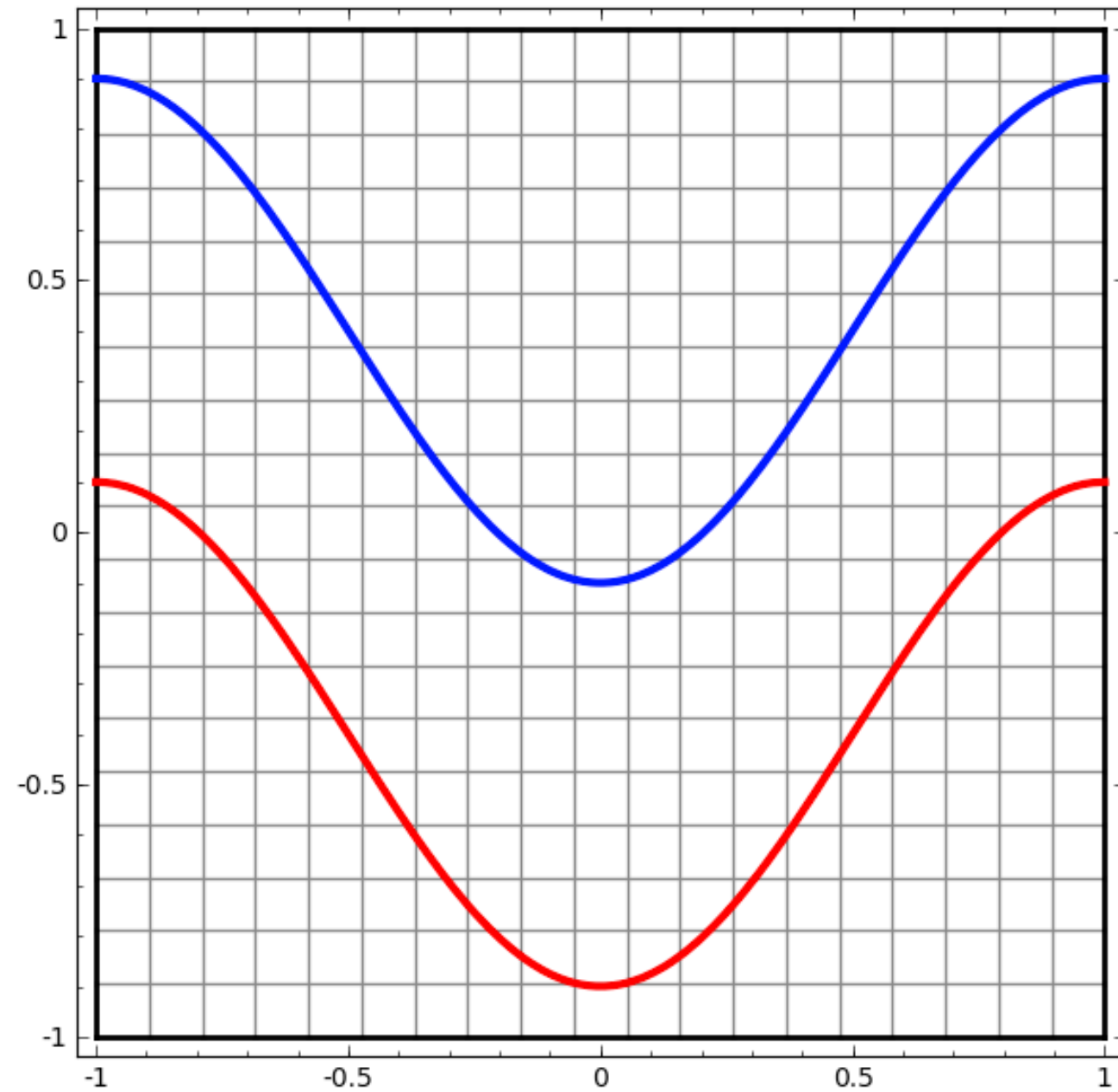
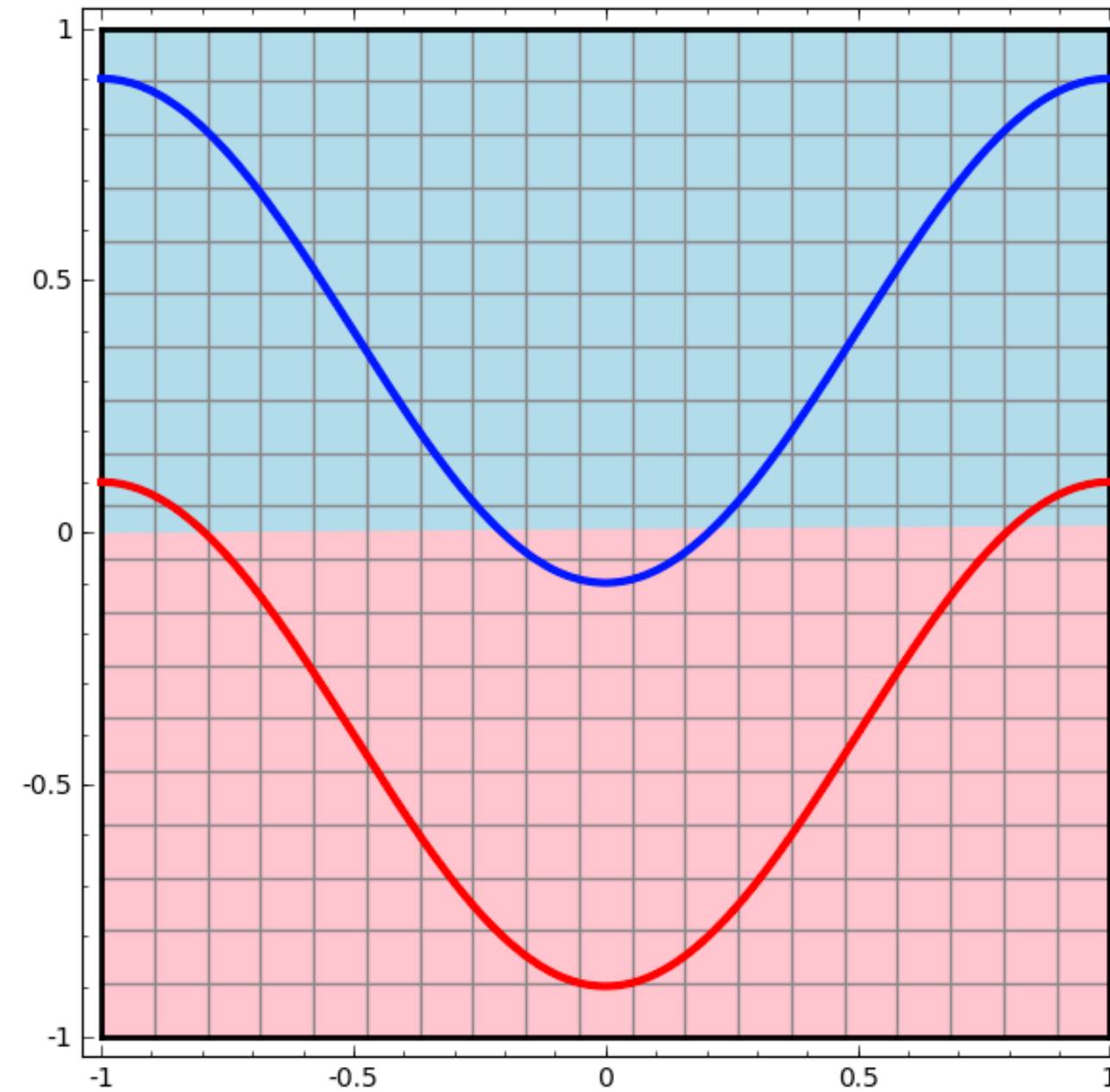depends on the
form of the loss

derivative of the
non-linearity

$$\frac{\partial}{\partial \mathbf{W}^{(L)}} (\mathbf{W}^{(L)\mathsf{T}} \mathbf{x}^{(L-1)})$$

$$= \mathbf{x}^{(L-1)\mathsf{T}}$$

note $\nabla_{\mathbf{W}^{(L)}} \mathcal{L} \equiv \dfrac{\partial \mathcal{L}}{\partial \mathbf{W}^{(L)}}$ is notational convention
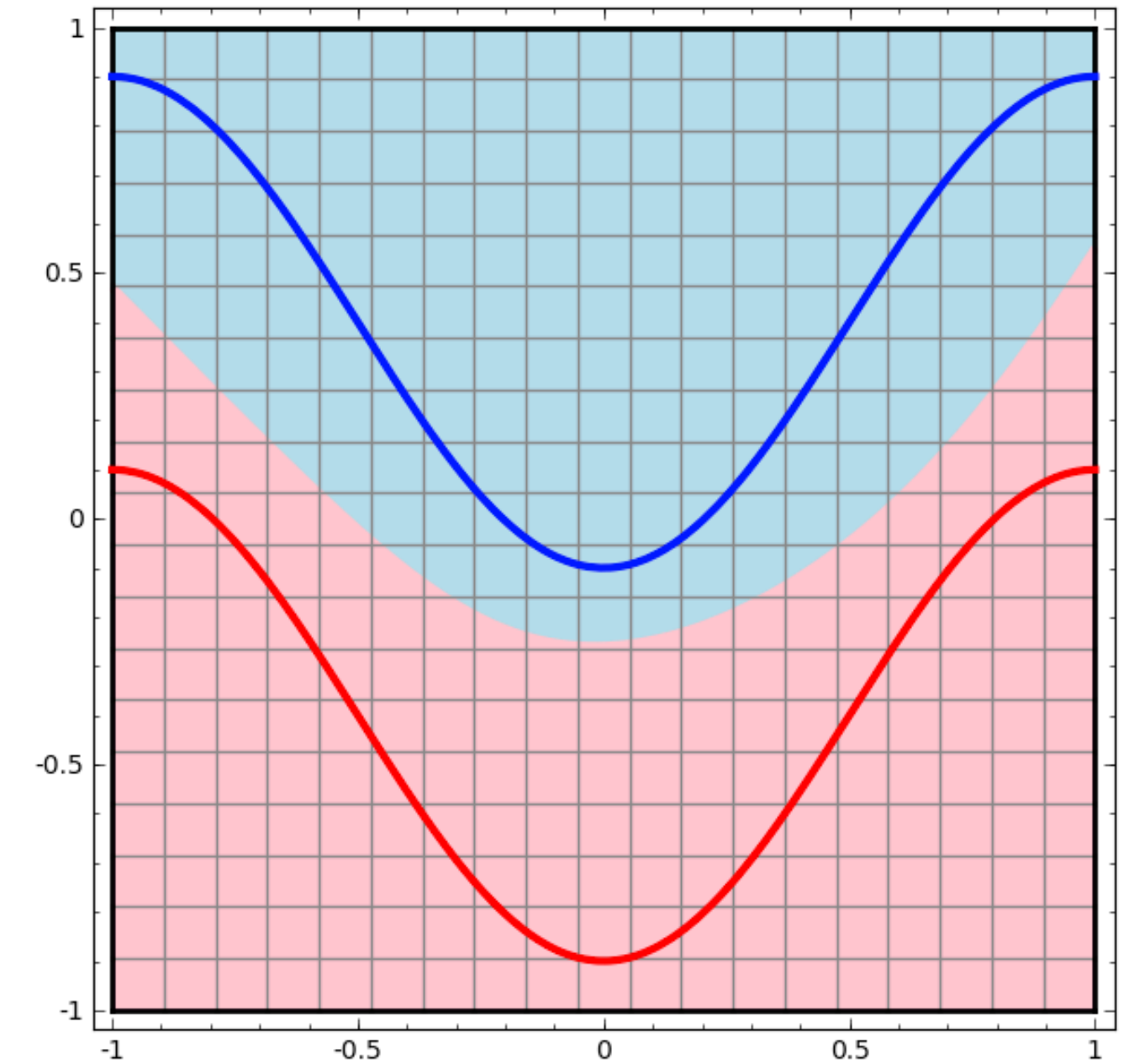
Data

Linear classifier

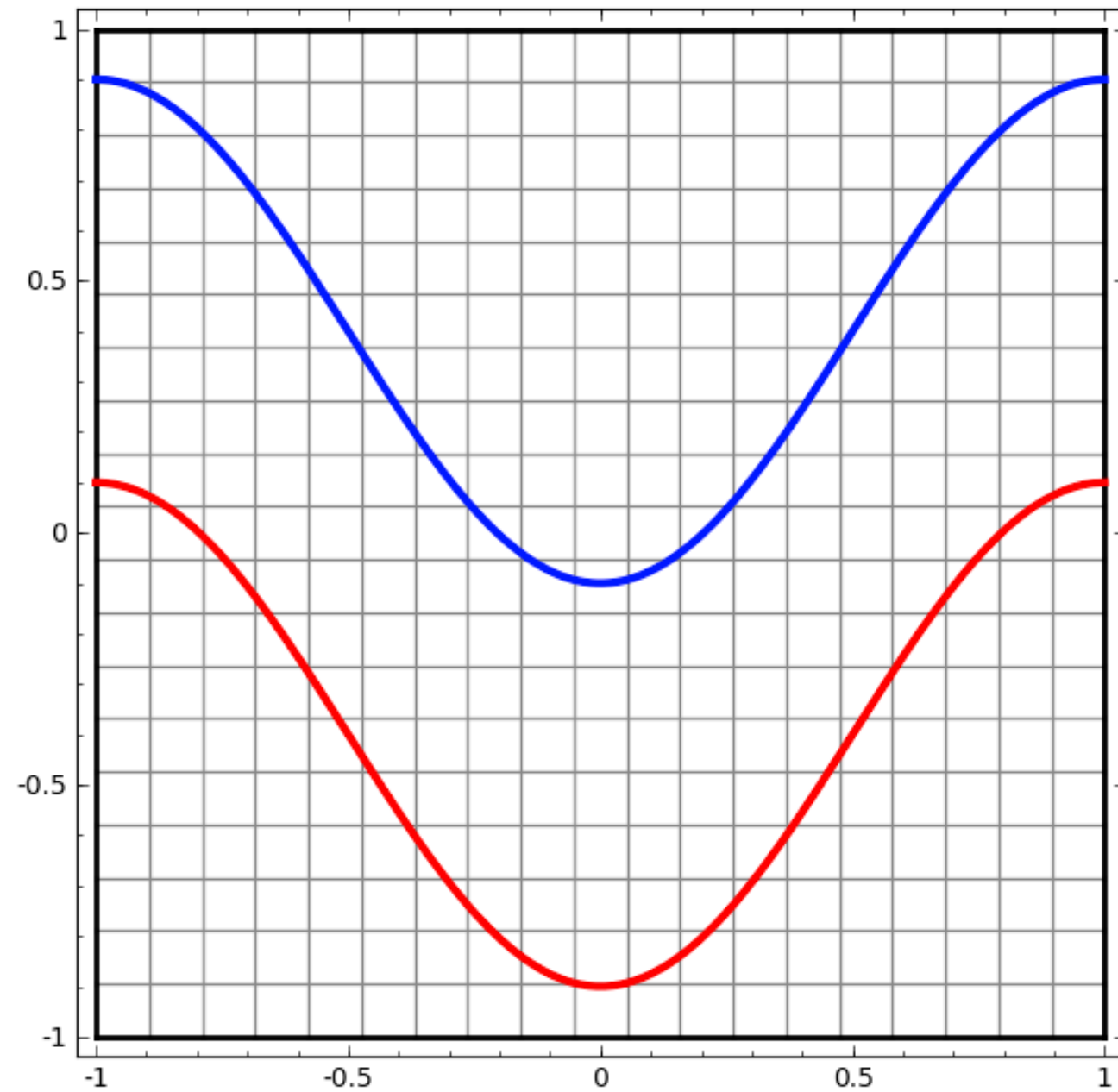Embedding + Linear classifier

$$y = \text{softmax}(w^\mathsf{T} x)$$

$$y = \text{softmax}(w^\mathsf{T} \phi(x))$$
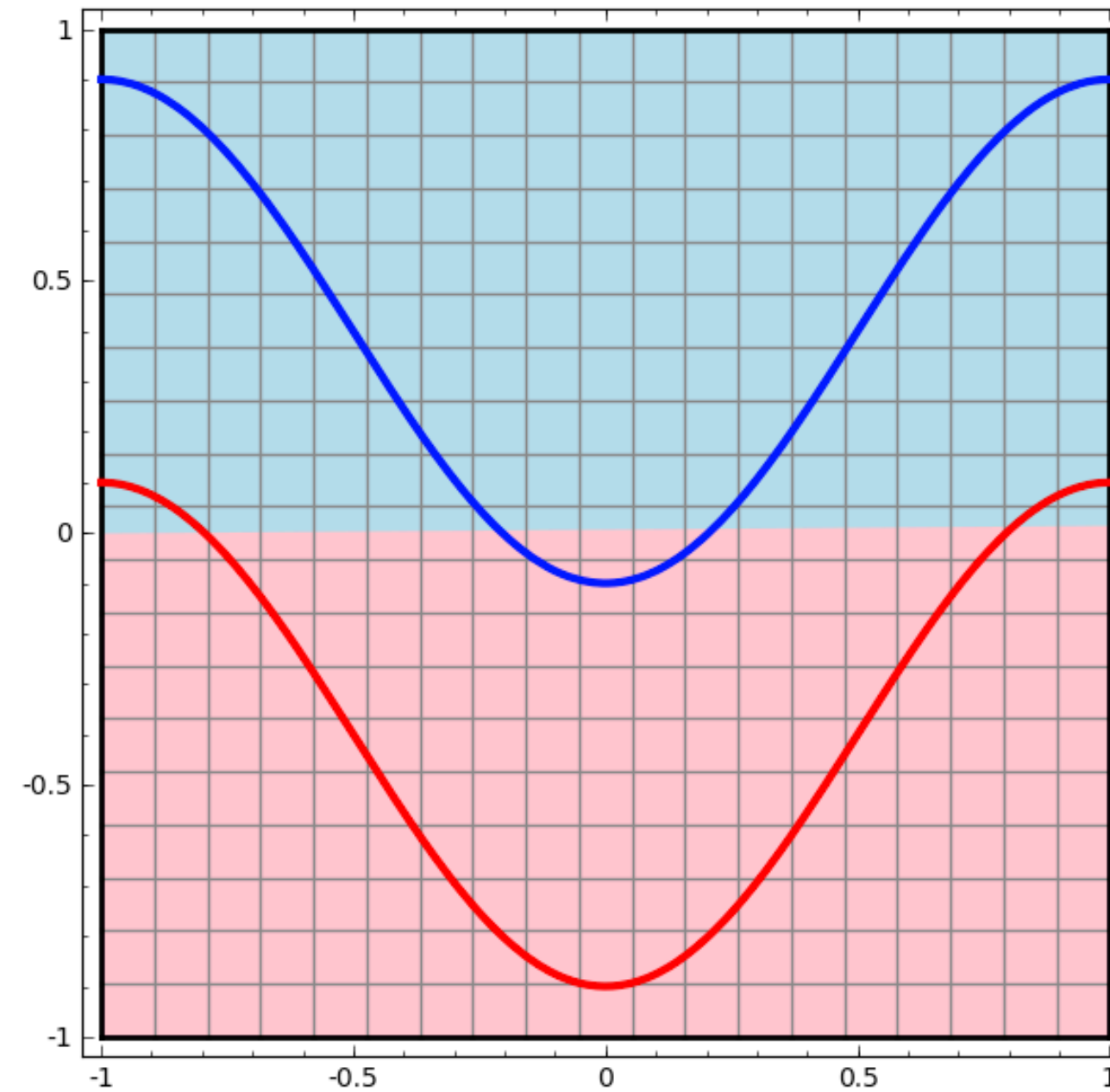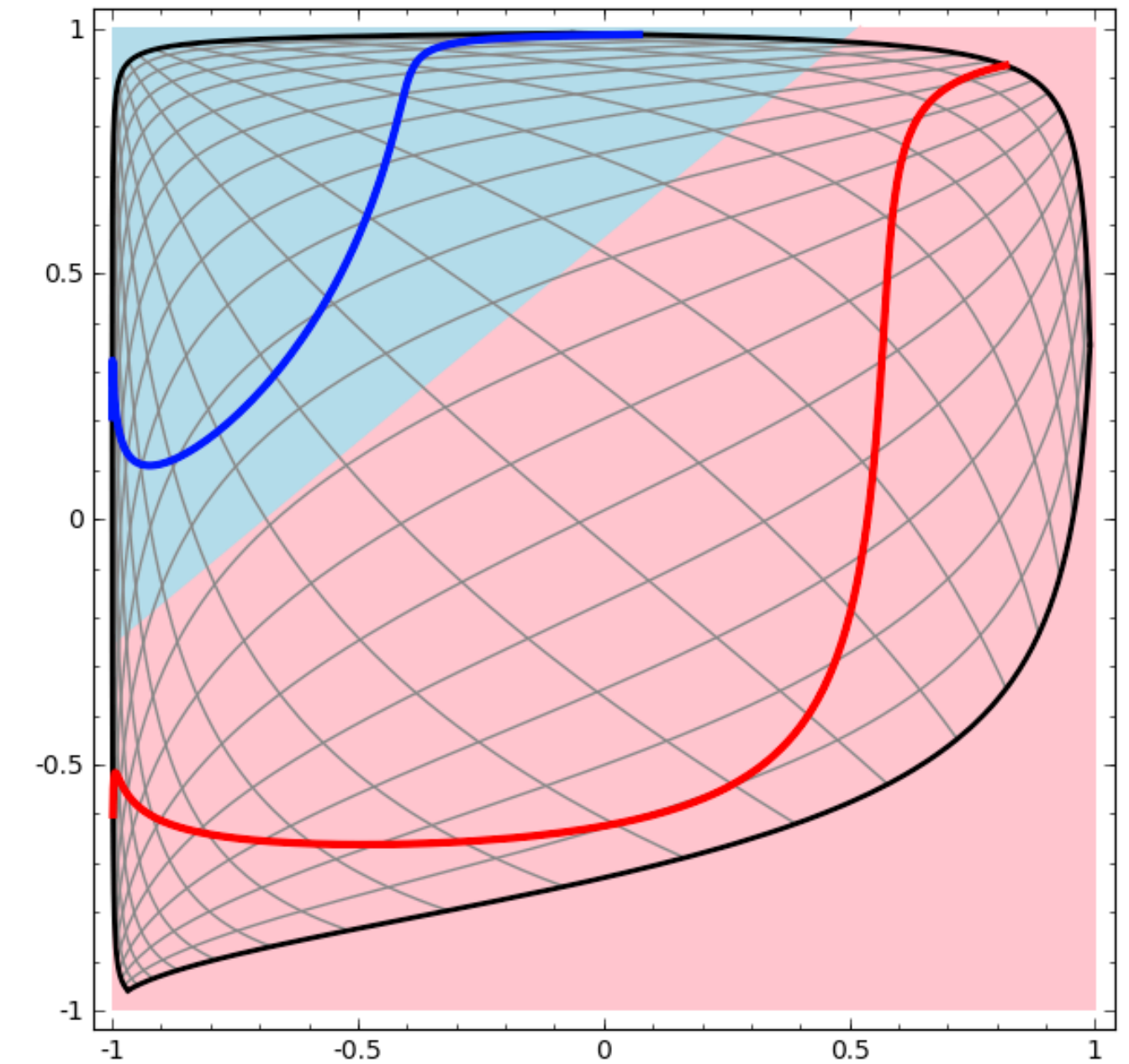
We have seen the polynomial embedding:

$$\phi(x) = (1, x, x^2, \ldots, x^n)$$

Data

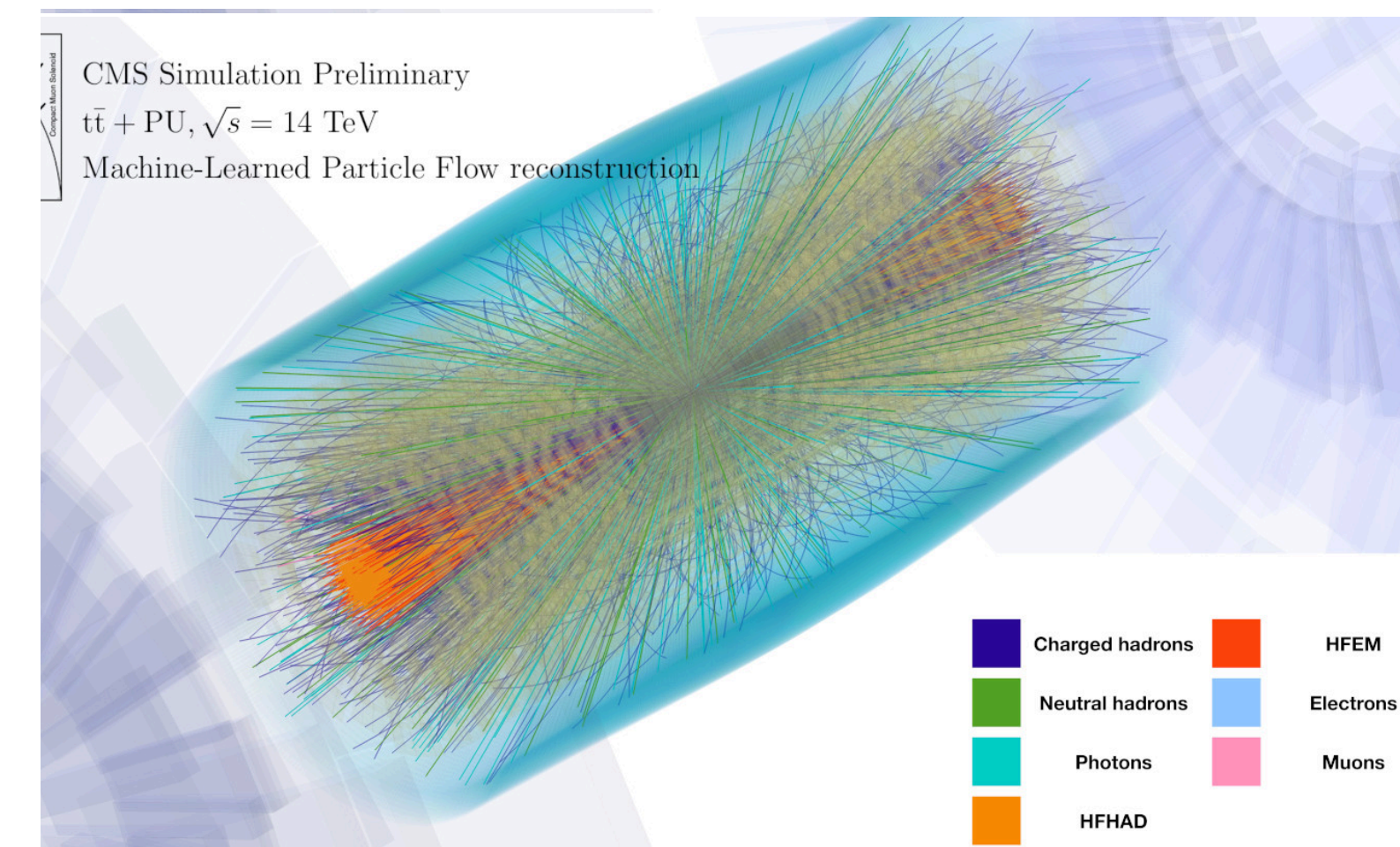Linear classifier

2-layer network

$$y = \text{softmax}(w^\mathsf{T}x)$$

$$y = \sigma(W_2\sigma(W_1x))$$

colah.github.io/posts/2014-03-NN-Manifolds-Topology

# NNS IN THE WILD

- ▸ B-jet energy regression [arXiv:1912.06046]

- ▸ Jet classification [arXiv:2004.08262]

- ▸ Jet mass regression [https://cds.cern.ch/record/2777006]

- ▸ Tracking

- ▸ Clustering

- ▸ Particle-flow reconstruction [arXiv:2203.00330]

- ▸ Anomaly detection

- ▸ Fast simulation

- ▸ Trigger applications

- ▸ Background modeling

▸ Jet images = pixelated versions of calorimeter hits in 2D ($\eta$, $\phi$)

▸ Much lower level

top jet

vs.

*u,d* or *s* jet

top jet (on average)

QCD jet (on average)

*Jet ETmiss*



Convolutions

Convolved
Feature Layers

Max-Pooling

Repeat

*W'→ WZ* event

**Locality** and **translation invariance** as *inductive biases*

## locality

*nearby areas tend to contain stronger patterns*

inputs can be restricted to regions

maintain spatial ordering

## translation invariance

*relative positions are relevant*

same filters can be applied throughout the input

same weights

▸ CNNs among the best performing algorithms

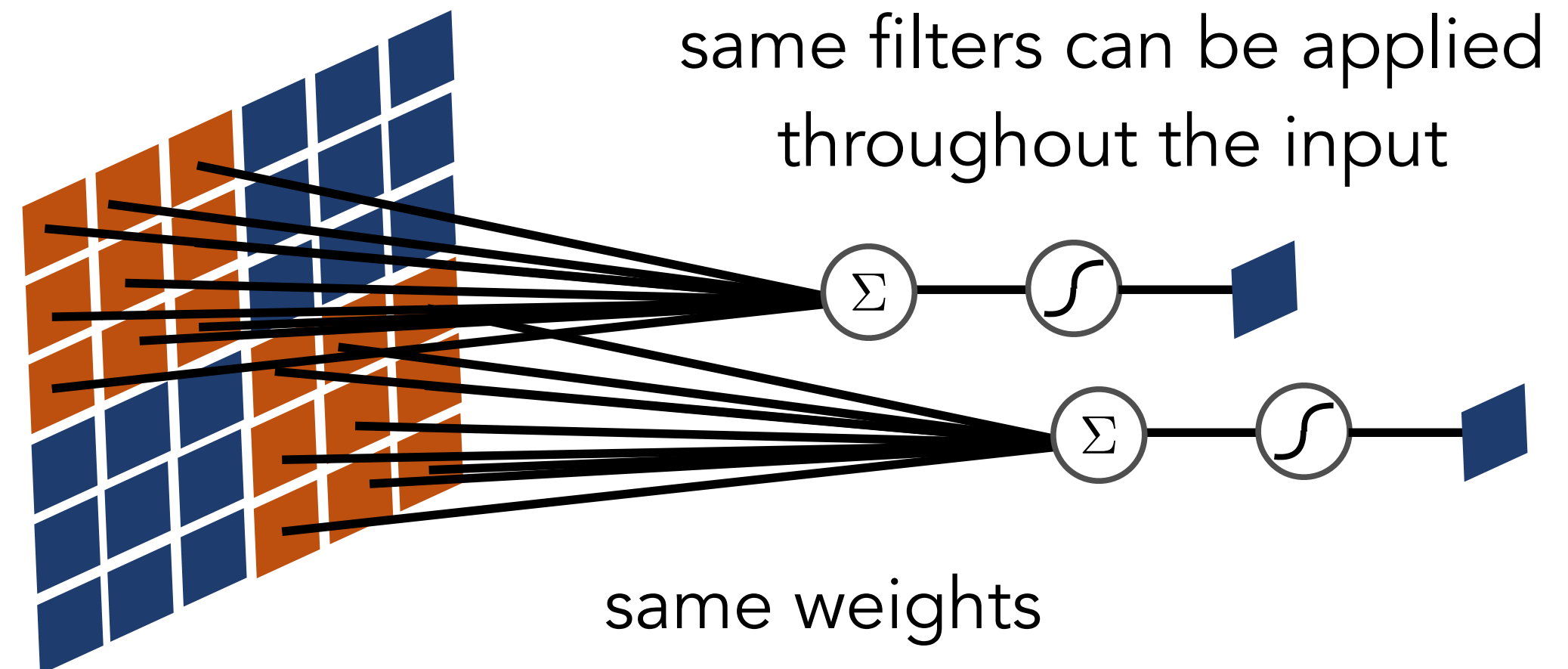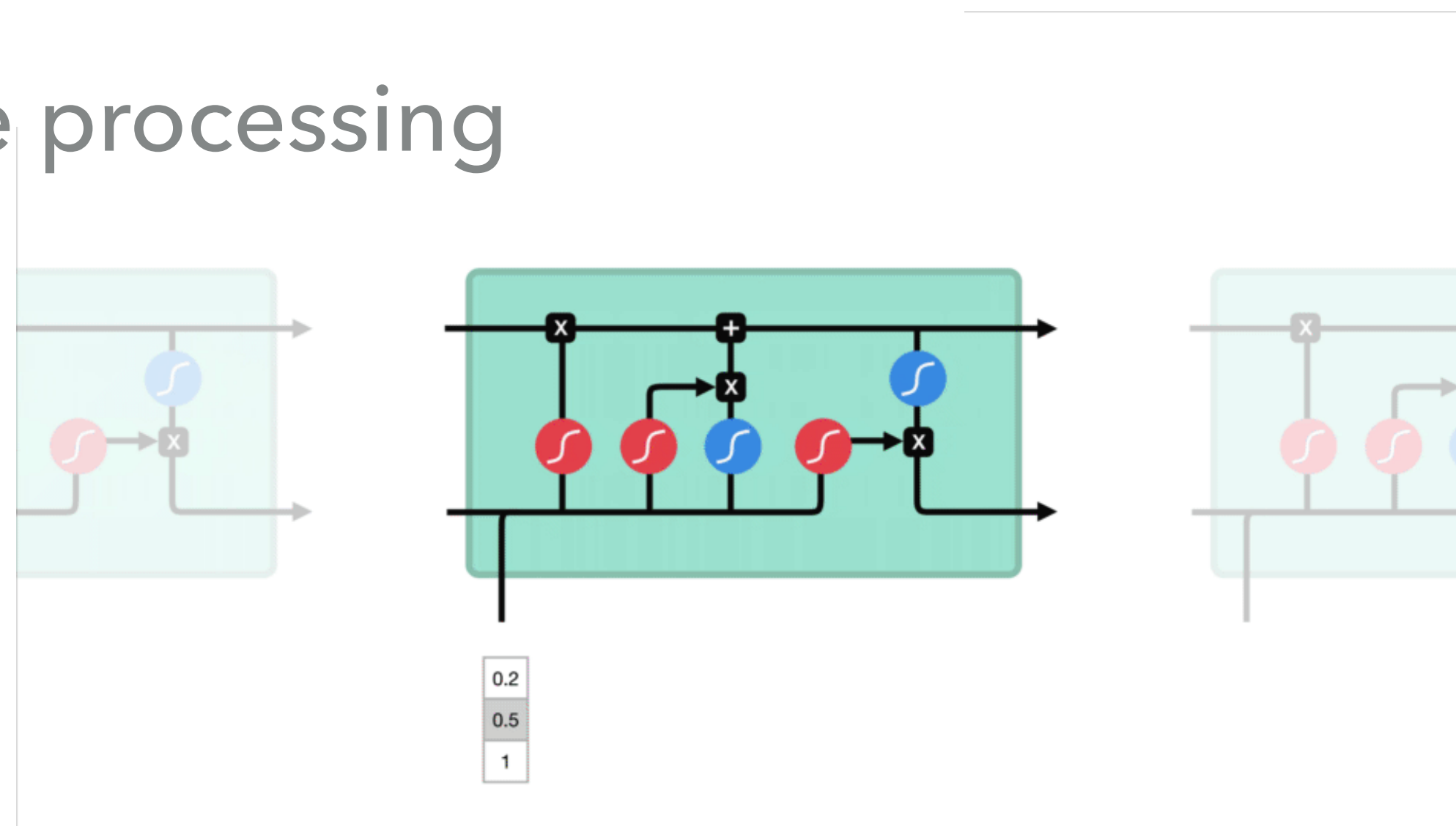| | AUC | Acc | $1/\epsilon_B$ ($\epsilon_S = 0.3$) | | | #Param |
|---|---|---|---|---|---|---|
| | | | single | mean | median | |
| CNN [16] | 0.981 | 0.930 | 914±14 | 995±15 | 975±18 | 610k |
| ResNeXt [31] | 0.984 | 0.936 | 1122±47 | 1270±28 | 1286±31 | 1.46M |
| TopoDNN [18] | 0.972 | 0.916 | 295±5 | 382± 5 | 378 ± 8 | 59k |
| Multi-body $N$-subjettiness 6 [24] | 0.979 | 0.922 | 792±18 | 798±12 | 808±13 | 57k |
| Multi-body $N$-subjettiness 8 [24] | 0.981 | 0.929 | 867±15 | 918±20 | 926±18 | 58k |
| TreeNiN [43] | 0.982 | 0.933 | 1025±11 | 1202±23 | 1188±24 | 34k |
| P-CNN | 0.980 | 0.930 | 732±24 | 845±13 | 834±14 | 348k |
| ParticleNet [47] | 0.985 | 0.938 | 1298±46 | 1412±45 | 1393±41 | 498k |
| LBN [19] | 0.981 | 0.931 | 836±17 | 859±67 | 966±20 | 705k |
| LoLa [22] | 0.980 | 0.929 | 722±17 | 768±11 | 765±11 | 127k |
| LDA [54] | 0.955 | 0.892 | 151±0.4 | 151.5±0.5 | 151.7±0.4 | 184k |
| Energy Flow Polynomials [21] | 0.980 | 0.932 | 384 | | | 1k |
| Energy Flow Network [23] | 0.979 | 0.927 | 633±31 | 729±13 | 726±11 | 82k |
| Particle Flow Network [23] | 0.982 | 0.932 | 891±18 | 1063±21 | 1052±29 | 82k |
| GoaT | 0.985 | 0.939 | 1368±140 | | 1549±208 | 35k |

▸ In deep learning, tailoring algorithms to the structure (and symmetries) of the data has led to groundbreaking performance
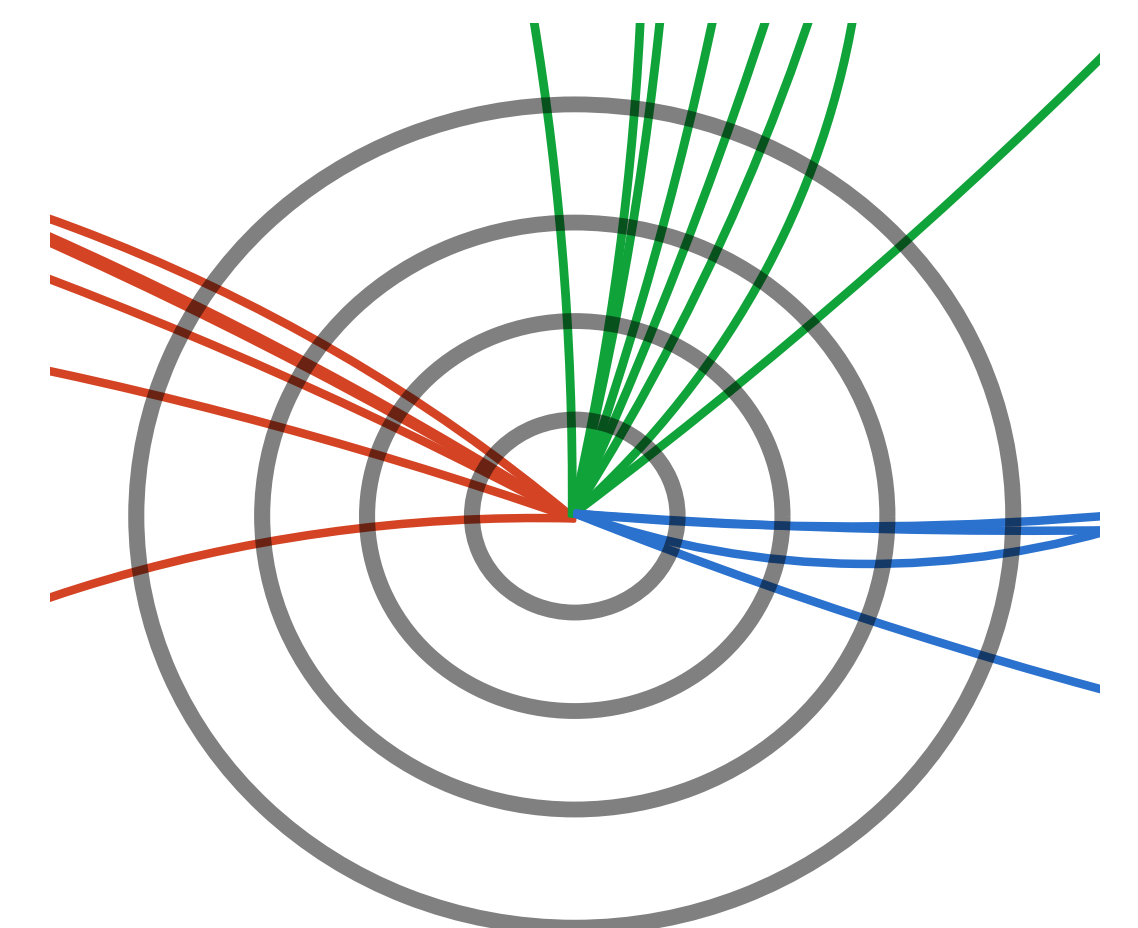
  ▸ CNNs for images

  ▸ RNNs for language processing

▸ What about high energy physics data like jets?

▸ Distributed unevenly in space
▸ Sparse
▸ Variable size
▸ No defined order
▸ Interconnections
  → Graphs

0.2
0.5
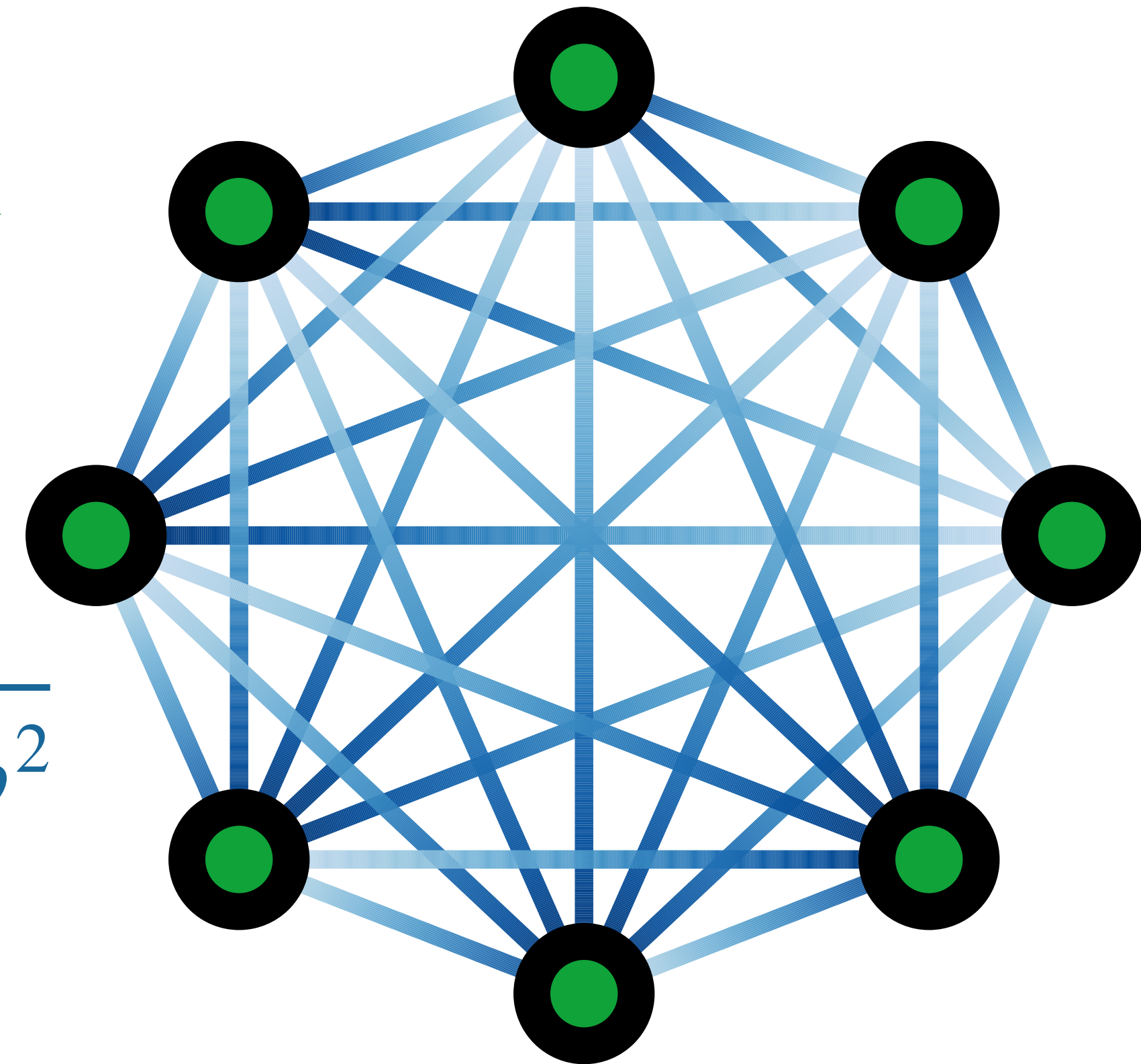1

▸ Node features $\mathbf{v}_i$: particle 4-momentum

$$p = [E, p_x, p_y, p_z] \equiv [p_{\mathrm{T}}, \eta, \phi, m]$$

▸ Edge features $\mathbf{e}_k$: pseudoangular distance between particles

$$\sqrt{\Delta\eta^2 + \Delta\phi^2}$$

▸ Graph (global) features $\mathbf{u}$: jet mass

$$m = \sqrt{\sum_{i \in \mathrm{jet}} E_i^2 - p_{x,i}^2 - p_{y,i}^2 - p_{z,i}^2}$$

▸ Features: triplet of global features, node features, and edge features: $(\mathbf{u}, V, E)$

▸ Graph connectivity: adjacency matrix

$A = \{a_{ij} = 1$ if $i$ is connected to $j\}$

  ▸ Sparse representation:

  "receiver" indices $r$ and "sender" indices $s$

  e.g. $k^{th}$ edge connects node $s_k$ to node $r_k$

Nodes
[ 0 , 1 , 1 , 0 , 0 , 1 , 1 , 1 ]

Edges
[ 2 , 1 , 1 , 2 , 2 , 1 , 1 ]

Adjacency List
[ [1, 0] , [2, 0] , [4, 3] , [6, 2] ,
  [7, 3] , [7, 4] , [7, 5] ]

Global
1

Attributes

▸ For all neighbors $j$ of node $i$ compute a "message" via a NN: $\phi(x_i, x_j)$

▸ Update the node features by summing all messages:

$$h_i = \sum_j \phi(x_i, x_j)$$



$\phi(x_i, x_j)$

$\mathbf{x}_i$

GNN

$\mathbf{h}_i$

Inputs
$(\mathbf{X}, \mathbf{A})$

Latents
$(\mathbf{H}, \mathbf{A})$

"message passing"

▸ Node-level tasks

  ▸ Identify "pileup" particles

▸ Graph-level tasks

  ▸ **Jet tagging**



**Node** classification
$$\mathbf{z}_i = f(\mathbf{h}_i)$$

**Graph** classification
$$\mathbf{z}_G = f\left(\bigoplus_{i \in \mathcal{V}} \mathbf{h}_i\right)$$

**Link** prediction
$$\mathbf{z}_{ij} = f(\mathbf{h}_i, \mathbf{h}_j, \mathbf{e}_{ij})$$

$\phi(x_i, x_j)$  $\mathbf{x}_i$

**GNN**

$\mathbf{h}_i$

**Inputs**
$$(\mathbf{X}, \mathbf{A})$$

**Latents**
$$(\mathbf{H}, \mathbf{A})$$

▸ Edge-level tasks

  ▸ Identify good track doublets

▸ ParticleNet, using "dynamic edge convolutions:" graph is constructed based on "closeness" in an abstract "latent" space

▸ Identifies H(bb) with an efficiency of ~50% while rejecting 99.9% of background

▸ Explainable AI (XAI) refers to the set of techniques employed to provide explanations for ML model predictions

▸ Layerwise relevance propagation (LRP) [1] computes relevance (R) scores for each neuron in a ML model

▸ Neuron's R score is a measure of its contribution to the model's output,
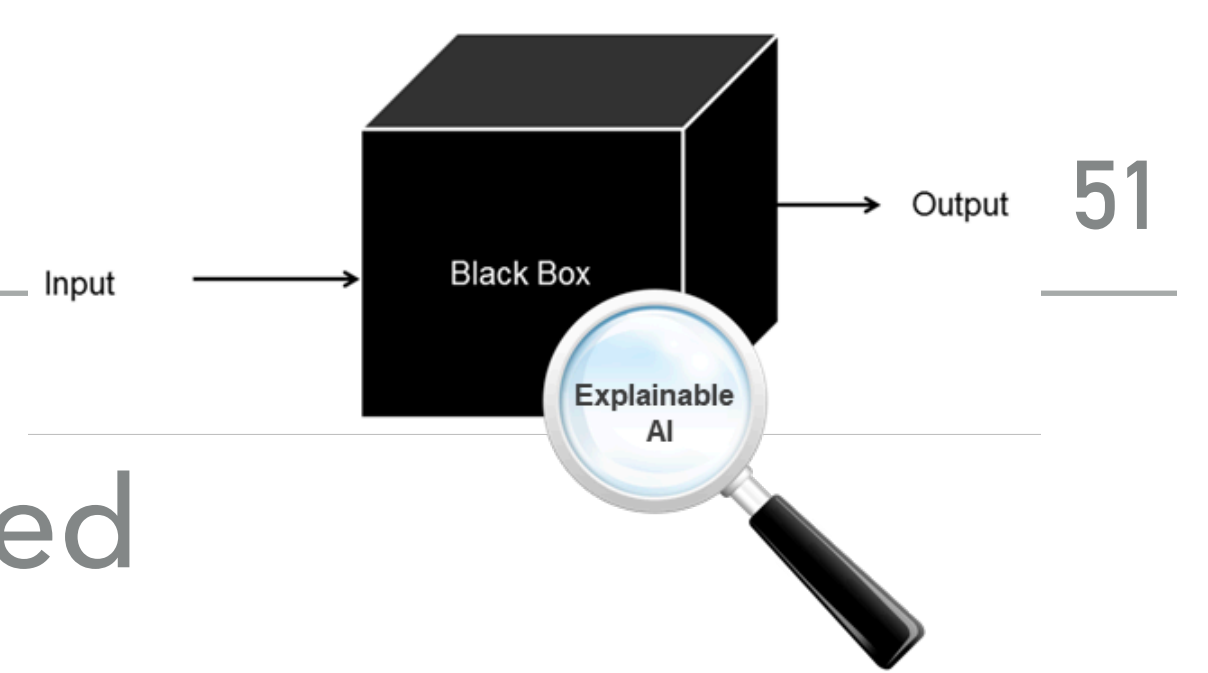
$$R_j^{(l)} = \sum_k \frac{z_{jk}}{\sum_m z_{mk}} R_k^{(l+1)} \text{ with } z_{jk} = x_j^{(l)} w_{jk}^{(l+1)}$$

▸ Flow of R scores for a multilayer perceptron (MLP)



[1] https://doi.org/10.1007/978-3-030-28954-6_10

▸ Is the model learning to connect particles from different ***subjets*** more often for top quark jets than for QCD jets?

▸ Use CA algorithm to decluster each jet into exactly 3 subjets

▸ For a top quark jet sample, relevant edges connect different subjets



▸ Trained model connects different subjets more often for top quarks than for QCD

▸ Symmetry-equivariant networks

　　▸ More economical (fewer, but more expressive parameters), interpretable, and trainable

**Invariance**

$$f(\rho_g(x)) = f(x)$$

**Equivariance**

$$f(\rho_g(x)) = \rho'_g(f(x))$$

▸ Lorentz-invariant networks:

  ▸ Boosting all particles into a new frame should give the same result

▸ Lorentz-equivariant networks:

  ▸ Boosting all particles into a new frame should give an output that transforms the same way

$$\begin{pmatrix} \gamma & -\gamma\beta & 0 & 0 \\ -\gamma\beta & \gamma & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}$$

Lorentz boost

arXiv:2201.08187

▸ State-of-the-art performance for top quark tagging

▸ Lorentz group invariance confirmed

I.    BASICS
II.   DATA REPRESENTATIONS & SYMMETRIES
III.  ANOMALY DETECTION
IV.   GENERATIVE MODELING
V.    SUMMARY & OUTLOOK

▸ Supervised = full label information

▸ Semi-supervised = partial labels

▸ Weakly-supervised = noisy labels

▸ Unsupervised = no labels

  ▸ Example: autoencoders compress data and then uncompress it

  ▸ Assumption: if $x$ is far from $\mathrm{Decoder}(\mathrm{Encoder}(x))$, then $x$ has low $p_{\mathrm{bkgd}}(x)$



background model independence

Some searches (train signal versus data)

**many new ideas!**

Most searches ("train" with simulations)

Train data versus background simulation

signal model independence

▸ Challenge with "black box" signals run in 2020–2021

▸ Plethora of new techniques

**3 Unsupervised**

3.1 Anomalous Jet Identification via Variational Recurrent Neural Network

3.2 Anomaly Detection with Density Estimation

3.3 BuHuLaSpa: Bump Hunting in Latent Space

3.4 GAN-AE and BumpHunter

3.5 Gaussianizing Iterative Slicing (GIS): Unsupervised In-distribution Anomaly Detection through Conditional Density Estimation

3.6 Latent Dirichlet Allocation

3.7 Particle Graph Autoencoders

3.8 Regularized Likelihoods

3.9 UCluster: Unsupervised Clustering

**4 Weakly Supervised**

4.1 CWoLa Hunting

4.2 CWoLa and Autoencoders: Comparing Weak- and Unsupervised methods for Resonant Anomaly Detection

4.3 Tag N' Train

4.4 Simulation Assisted Likelihood-free Anomaly Detection

4.5 Simulation-Assisted Decorrelation for Resonant Anomaly Detection

**5 (Semi)-Supervised**

5.1 Deep Ensemble Anomaly Detection

5.2 Factorized Topic Modeling

5.3 QUAK: Quasi-Anomalous Knowledge for Anomaly Detection

5.4 Simple Supervised learning with LSTM layers

▸ Challenge: if new physics has an unexpected signature that doesn't align with existing triggers, precious BSM events may be discarded at trigger level

▸ Can we use unsupervised algorithms to detect non-SM-like anomalies?

  ▸ Autoencoders (AEs): compress input to a smaller dimensional latent space then decompress and calculate difference

  ▸ Variational autoencoders (VAEs): model the latent space as a probability distribution; possible to detect anomalies purely with latent space variables

Key observation: Can build an anomaly score from the latent space of VAE directly! No need to run decoder!

Encoder

μ

σ

$$R_z = \sum_i \frac{\mu_i^2}{\sigma_i^2}$$

▸ **hls4ml** for scientists or ML experts to translate ML algorithms into RTL firmware

**Machine learning model optimization, compression**

▸ CNNs as the basis for (V)AEs for anomaly detection

▸ Good anomaly detection performance for unseen signals
(LQ → bτ, A → 4l, **h± → τν**, h⁰ → ττ)

▸ **VAE** fits in latency and resource requirements for HL-LHC!



Block 1:
Conv2d (16,(3,3))
ReLU
AvPooling (3,1)

Block 2:
Conv2d 1 (32,(3,1))
ReLU
AvPooling (3,1)
Flatten (64)

Block 3:
Dense (8)
Dense 1 (64)
ReLU
Reshape (2,1,32)

Block 4:
Conv2d 2 (32,(3,1))
ReLU
UpSampling (3,1)
ZeroPad (0,0),(1,1)

Block 5:
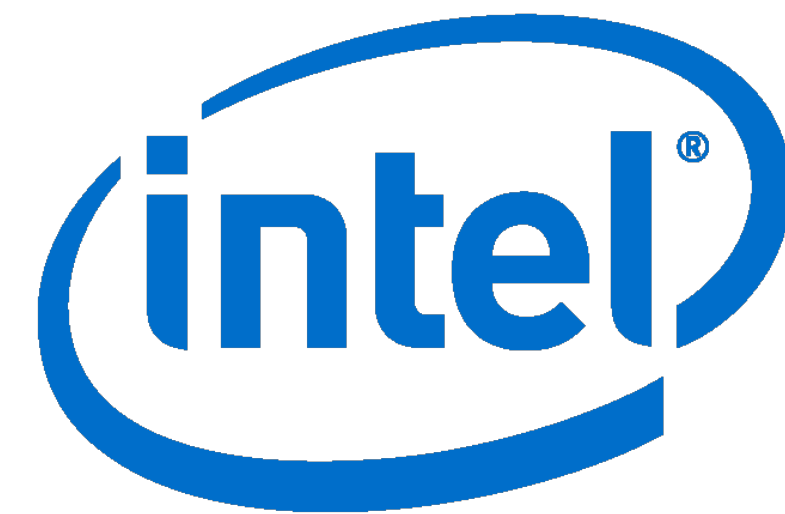Conv2d 3 (16,(3,1))
ReLU
UpSampling (3,1)
ZeroPad (1,0),(0,0)

Block 0:
Input 19x3x1
ZeroPadding (1,0)
BatchNorm

Output:
Conv2d 4 (1,(3,3))

CNN ROC $h^± → τν$

- IO VAE (AUC = 95%)
- VAE $D_{KL}$ (AUC = 86%)
- VAE $R_z$ (AUC = 86%)
- IO AE (AUC = 96%)

| Model | DSP [%] | LUT [%] | FF [%] | BRAM [%] | Latency [ns] | II [ns] | AUC [%] | TPR @ FPR=$10^{-5}$ |
|---|---|---|---|---|---|---|---|---|
| CNN VAE $R_z$ | 10 | **12** | **4** | **2** | **365** | **115** | 86 | 0.06% |
| CNN AE | **7** | 47 | 5 | 6 | 1480 | 895 | **96** | **0.10%** |

I. BASICS
II. DATA REPRESENTATIONS & SYMMETRIES
III. ANOMALY DETECTION
IV. GENERATIVE MODELING
V. SUMMARY & OUTLOOK

- ▸ Several different strategies:
  - ▸ Replace (part of) FullSim: increase speed, preserve accuracy
  - ▸ Replace (part of) FastSim: maintain speed, increase accuracy
  - ▸ Conditional: map generated → reconstructed events
  - ▸ **End-to-end: map random noise → reconstructed events directly**



| | | END–TO–END | | | | |
|---|---|---|---|---|---|---|
| | | | CONDITIONAL | | | |
| **Full Simulation** | | FULL DETECTOR SIM W/ ML | DIGITIZATION EMULATION | RECONSTRUCTION | | |
| HARD PROCESS GENERATION | SHOWERING/ HADRONIZATION/ UNDERLYING EVT | APPROXIMATE DETECTOR SIM W/ ML | PARTIAL DIGITIZATION EMULATION | SIMPLIFIED RECONSTRUCTION | | ANALYSIS/ NTUPLING |
| | Delphes | PARAMETRIZED SMEARING | | | | |

# GENERATIVE ADVERSARIAL NETWORKS

Note: failure modes!

thispersond...

- ▸ Train two neural networks in tandem:
  - ▸ one to generate realistic "fake" data
  - ▸ the other to discriminate "real" from "fake" data
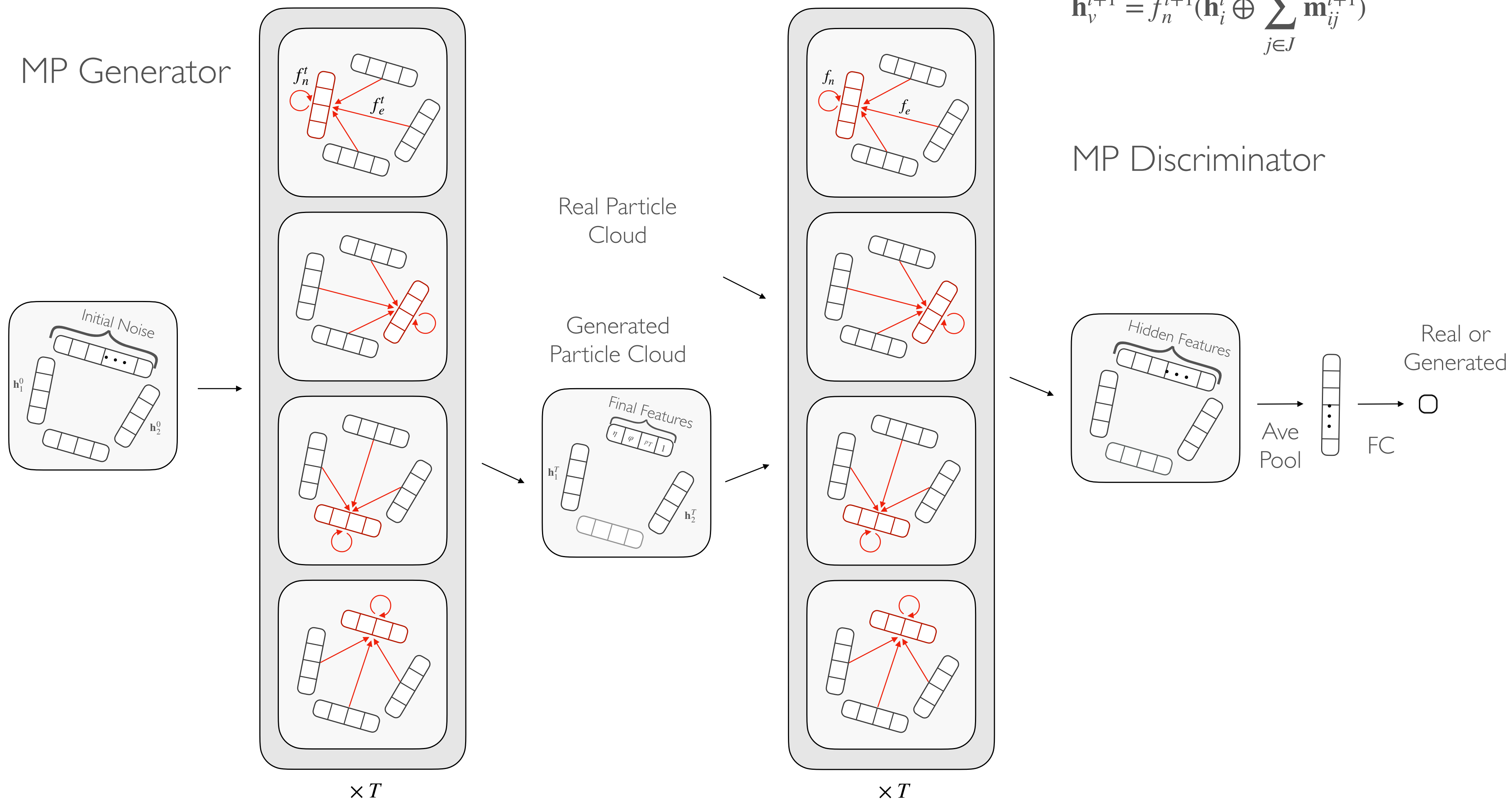
# GENERATIVE AI EVALUATION METRICS

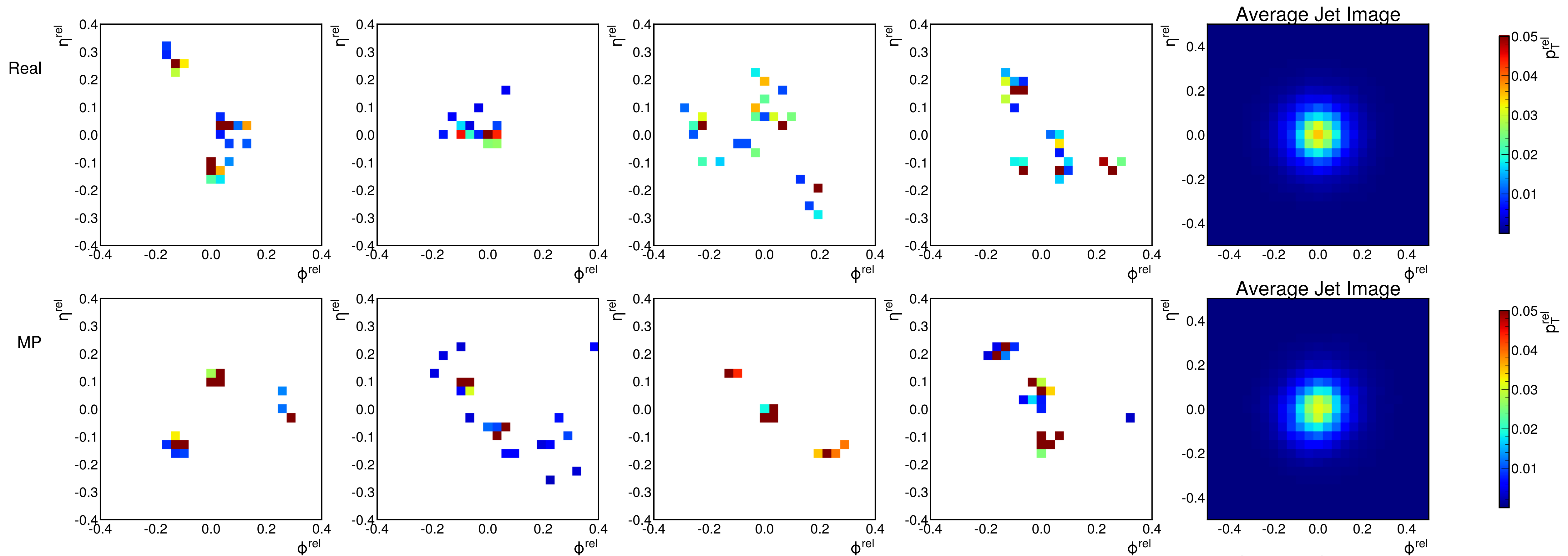▸ Evaluation of generative models is in general difficult

▸ We want to evaluate quantitatively:

  ▸ the **quality** of the data

  ▸ the **diversity** of the data

  ▸ ultimately, **physics performance**

| | Minimum Matching Distance | Coverage | Fréchet ParticleNet Distance | 1-Wassersstein Distance ($W_1$) |
|---|---|---|---|---|
| Quality | ✅ | | ✅ | ✅ |
| Diversity | | ✅ | ✅ | ✅ |
| Physics Perf. | | | | ✅ |

▸ As an alternative to voxelization, a graph-based GAN can be used to generate jets as particle clouds

$$\mathbf{m}_{ij}^{t+1} = f_e^{t+1}(\mathbf{h}_i^t \oplus \mathbf{h}_j^t)$$

$$\mathbf{h}_v^{t+1} = f_n^{t+1}(\mathbf{h}_i^t \oplus \sum_{j\in J} \mathbf{m}_{ij}^{t+1})$$

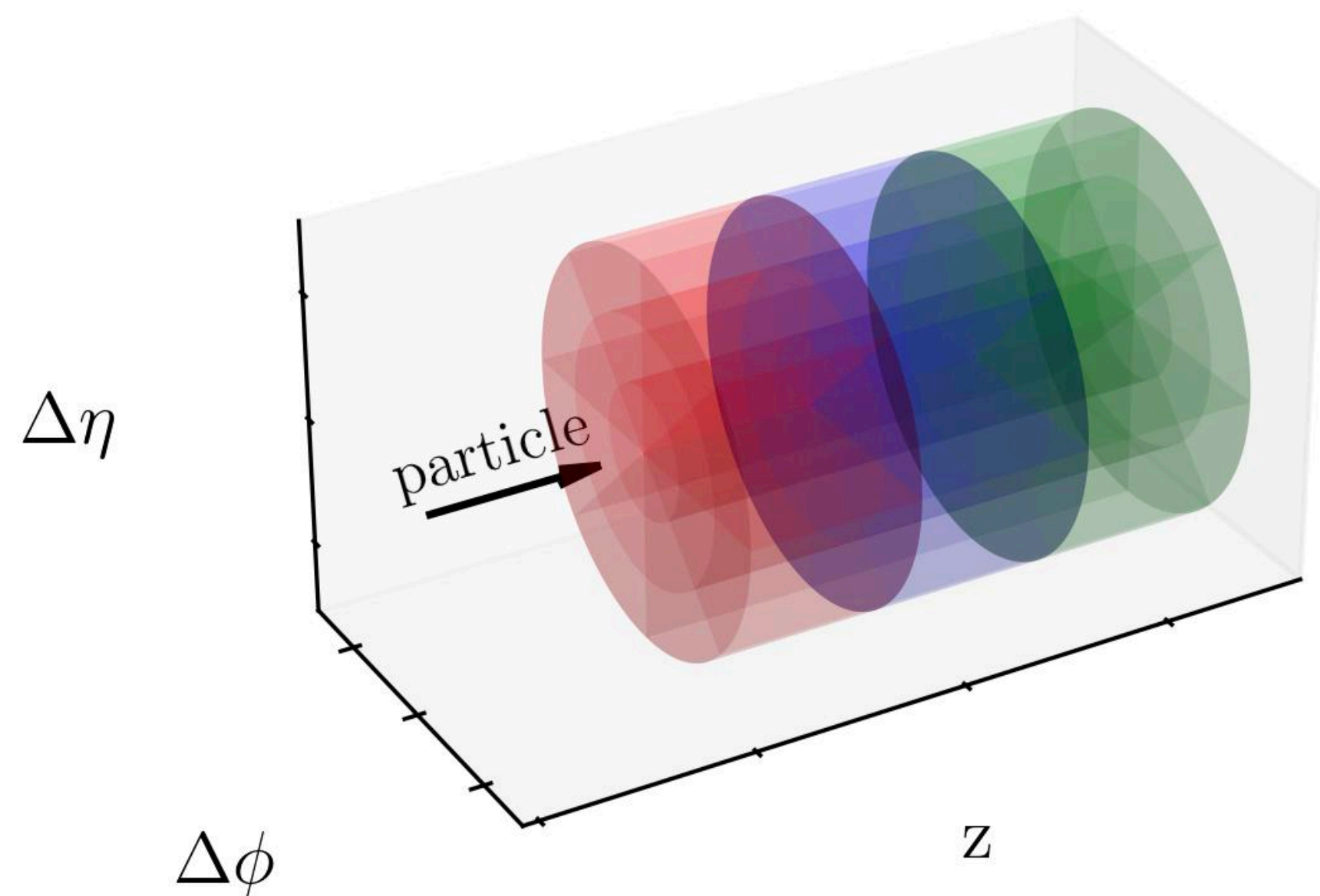▸ To easily visualize the generated particle clouds, we can make "jet images"

▸ Reproduces nontrivial properties like top quark jet mass and energy-flow polynomials ($t \rightarrow Wb \rightarrow qqb$ so 3 subjets + not always fully merged)



| Generator | Discriminator | $W_1$-P ($10^{-3}$) | $W_1$-M ($10^{-3}$) | $W_1$-EFP ($10^{-5}$) | FPND | Coverage | MMD | Time |
|---|---|---|---|---|---|---|---|---|
| | **Real** | **0.55 ± 0.07** | **0.51 ± 0.07** | **1.1 ± 0.1** | – | – | – | **O(1 s) per jet** |
| FC | PointNet | **1.5 ± 0.2** | 2.6 ± 0.1 | 8 ± 3 | 225 | 0.56 | 0.076 | |
| GraphCNN | PointNet | 37 ± 3 | 10.8 ± 0.5 | 39 ± 18 | 2M | 0.39 | 0.084 | |
| **MP** | **MP** | 2.1 ± 0.2 | **0.7 ± 0.1** | **1.8 ± 0.8** | **6.4** | 0.56 | **0.071** | **36 μs per jet** |
| MP | PointNet | **1.5 ± 0.1** | 1.0 ± 0.3 | 5 ± 2 | 12 | **0.58** | **0.071** | |

# CALO CHALLENGE

▸ Ongoing challenge for generative modeling of calorimeter showers in HEP!

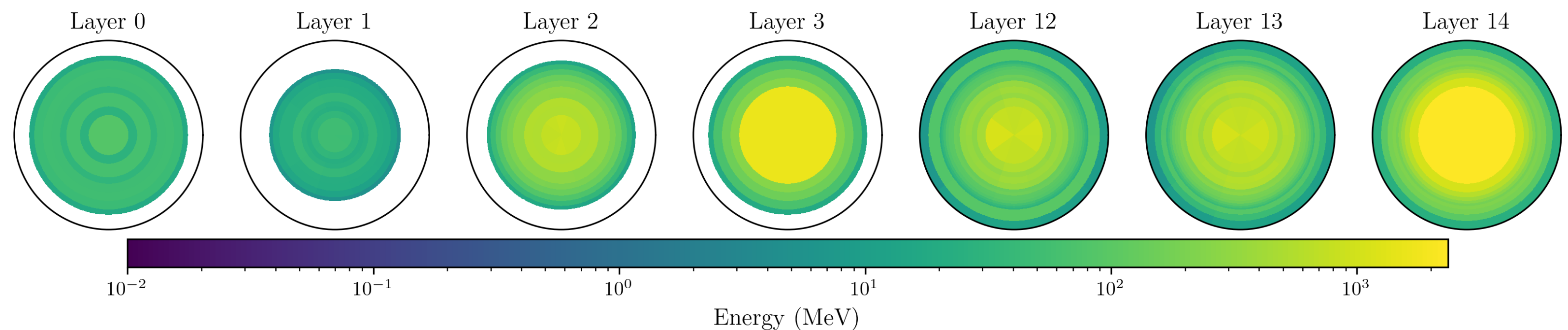▸ Many new approaches presented at CaloChallenge Workshop: https://agenda.infn.it/event/34036/



Shower average GEANT4 photon reference dataset

Shower average GEANT4 pion reference dataset

- Diffusion models have recently dethroned GANs for natural images

- Generative model is trained using a diffusion process that slowly perturbs the data by adding noise — model learns to **denoise**

- Generation of new samples by reversing the diffusion process

**Forward diffusion (training)**



t=0   t=0.25   t=0.75   t=1

**Reverse-time diffusion (data generation)**

- Distribution of deposited energies for generated particle energies (top) and the energy deposition in a single layer of a calorimeter (bottom) vs time step

# I. BASICS
# II. DATA REPRESENTATIONS & SYMMETRIES
# III. ANOMALY DETECTION
# IV. GENERATIVE MODELING
# V. SUMMARY & OUTLOOK

▸ Different representations of HEP data, from tabular data, image data, set data, graph data, paired with correspondir algorithms can achieve excellent performance

▸ Plethora of ML techniques in HEP from anomaly detection to generative modeling have exploded in recent years

▸ Availability of public datasets and challenges have advanced the state-of-the-art

## ts as Images

▸ Fast ML can accelerate science allowing us to test hypotheses faster, enhance performance of detectors/ accelerators, and save potentially overlooked data

φ) to a rectangular grid that allows for an image-ergy from particles are deposited in pixels in (η, φ) em as the pixel intensities in a greyscale analogue. st introduced by our group [JHEP 02 (2015) 118], even replace current s event reconstruction and computer vision.. We simulators he jet-axis, and normalize each image, as is often scriminative difference in pixel intensities.

Below, we have signal and backg difference-visualizat helps

B. Nachman (SLAC)