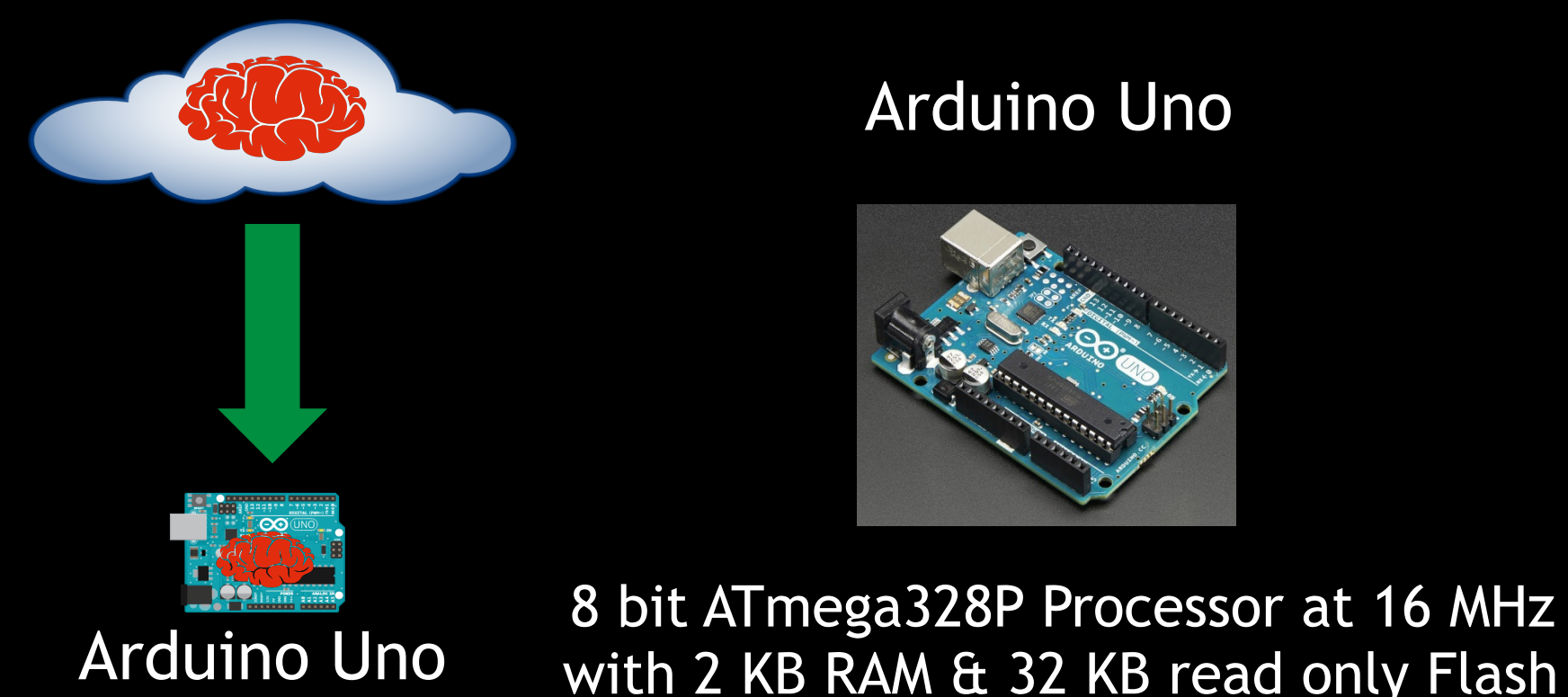# Efficient ML in 2 KB RAM for the Internet of Things
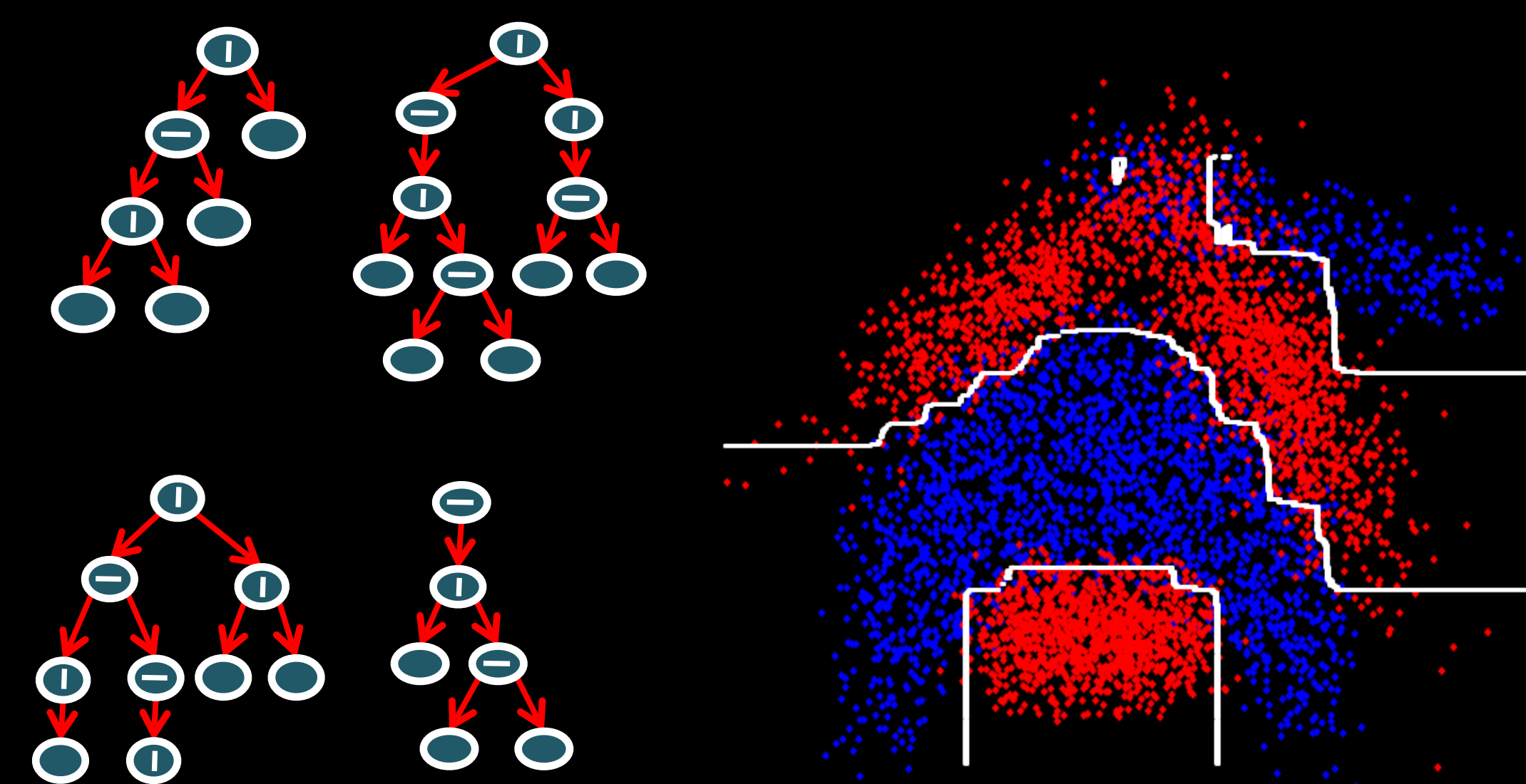
Ashish Kumar (MSR)
Saurabh Goyal (IIT Delhi)
Manik Varma (MSR)

## Objective

- To build an efficient tree classifier
  - Which can be trained on the cloud
  - But which can make predictions on tiny IoT devices

Arduino Uno

Arduino Uno

8 bit ATmega328P Processor at 16 MHz with 2 KB RAM & 32 KB read only Flash
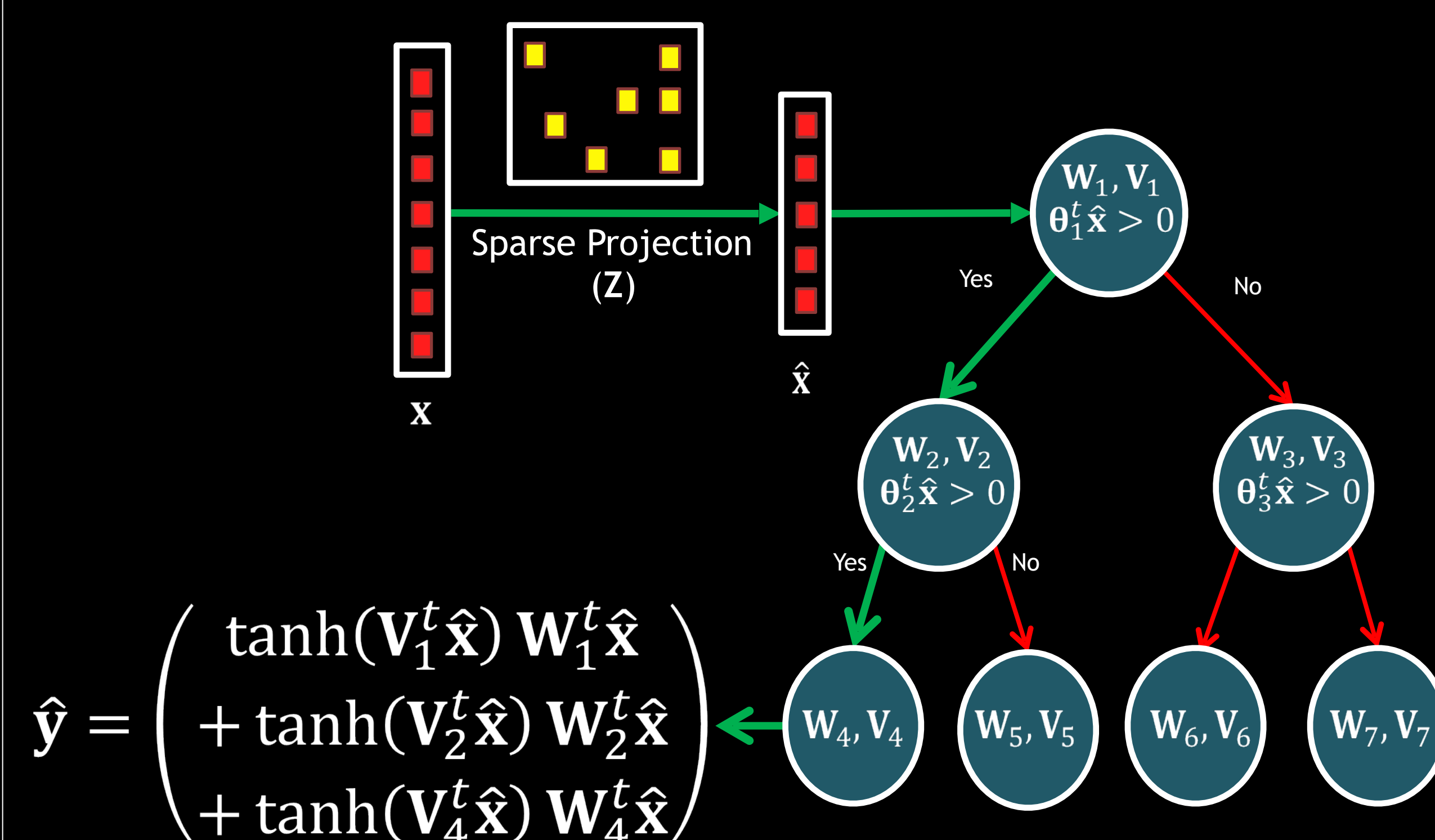
## Disadvantages of Tree Ensembles

- Tree Ensembles might not fit in Kilobytes and might not be accurate
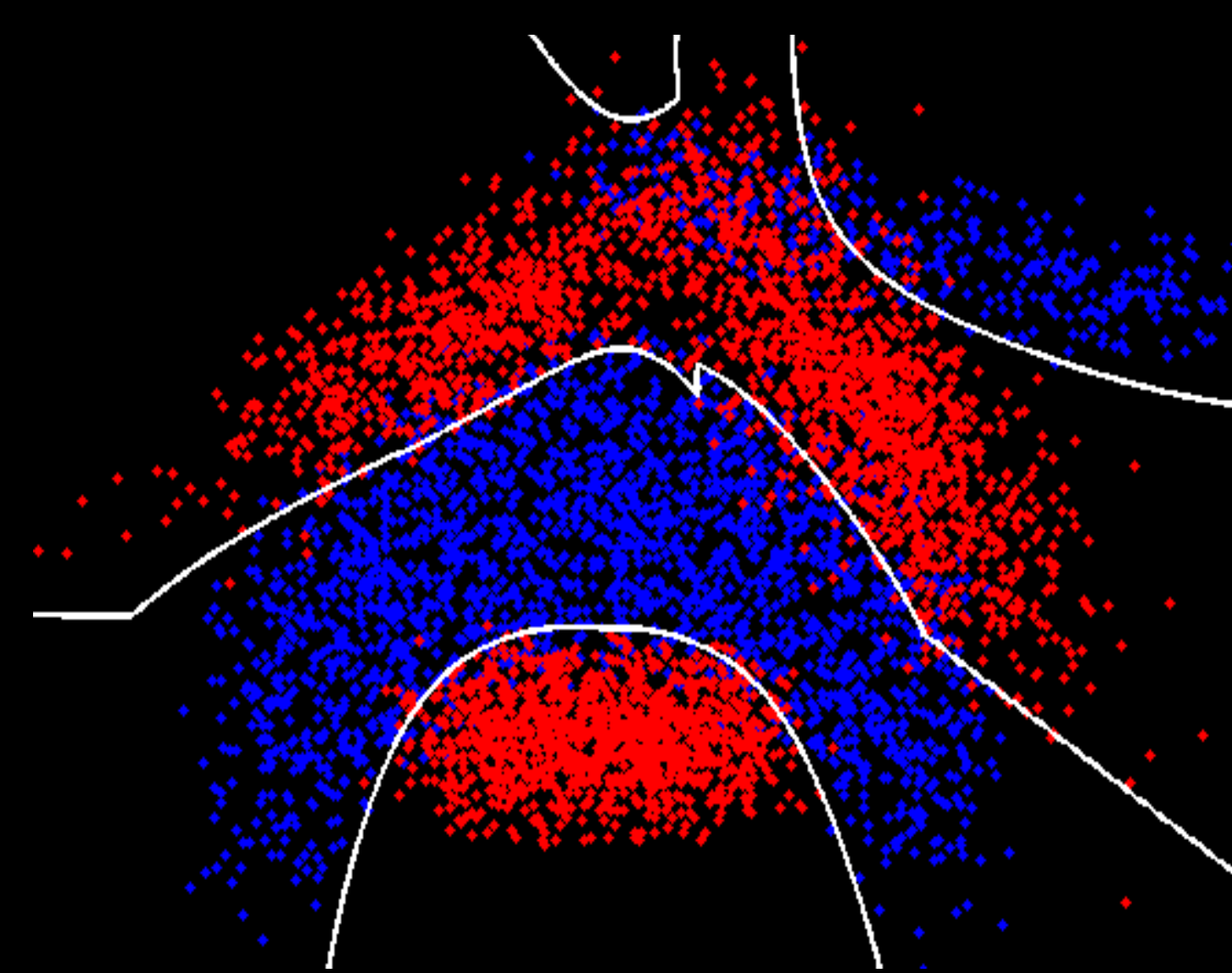
## Bonsai – Key Ideas

- We design Bonsai to be a single, shallow, sparse tree with powerful nodes for accurate prediction

- We reduce model size by learning Bonsai in a low-dimensional space into which all data is projected

- We jointly learn tree and projection parameters so as to maximize accuracy within the given budget

## Bonsai – A Compact Tree Model

x

Sparse Projection (Z)

$\hat{x}$

$W_1, V_1$ $\theta_1^t \hat{x} > 0$

Yes    No

$W_2, V_2$ $\theta_2^t \hat{x} > 0$

$W_3, V_3$ $\theta_3^t \hat{x} > 0$

Yes    No

$W_4, v_4$    $W_5, v_5$    $W_6, v_6$    $W_7, v_7$

$$\hat{y} = \begin{pmatrix} \tanh(V_1^t \hat{x}) \; W_1^t \hat{x} \\ + \tanh(V_2^t \hat{x}) \; W_2^t \hat{x} \\ + \tanh(V_4^t \hat{x}) \; W_4^t \hat{x} \end{pmatrix}$$

## Bonsai's Decision Boundaries

## The Bonsai Objective Function

$$\min_{\Theta, Z} \; P = \frac{1}{2} Tr(\Lambda \Theta^\top \Theta) + \frac{\lambda}{2} Tr(Z Z^\top) + \frac{1}{N} \sum_{i=1}^{N} \mathcal{L}(x_i, y_i, \hat{y}_i; \Theta, Z)$$

$$s.t. \quad \left\| \Theta \right\|_0 < B, \; \left\| Z \right\|_0 < S, \quad \Theta = [W, V, \theta]$$

- $\mathcal{L}$ is the loss function for classification, regression and ranking which can be optimized via SGD

- We place explicit budget constraints on the tree parameters $\Theta$ and sparse projection matrix $Z$

## Bonsai Optimization

Algorithm – Repeat the following steps till convergence

1. Mini batch gradient descent: $K$ steps with fixed support
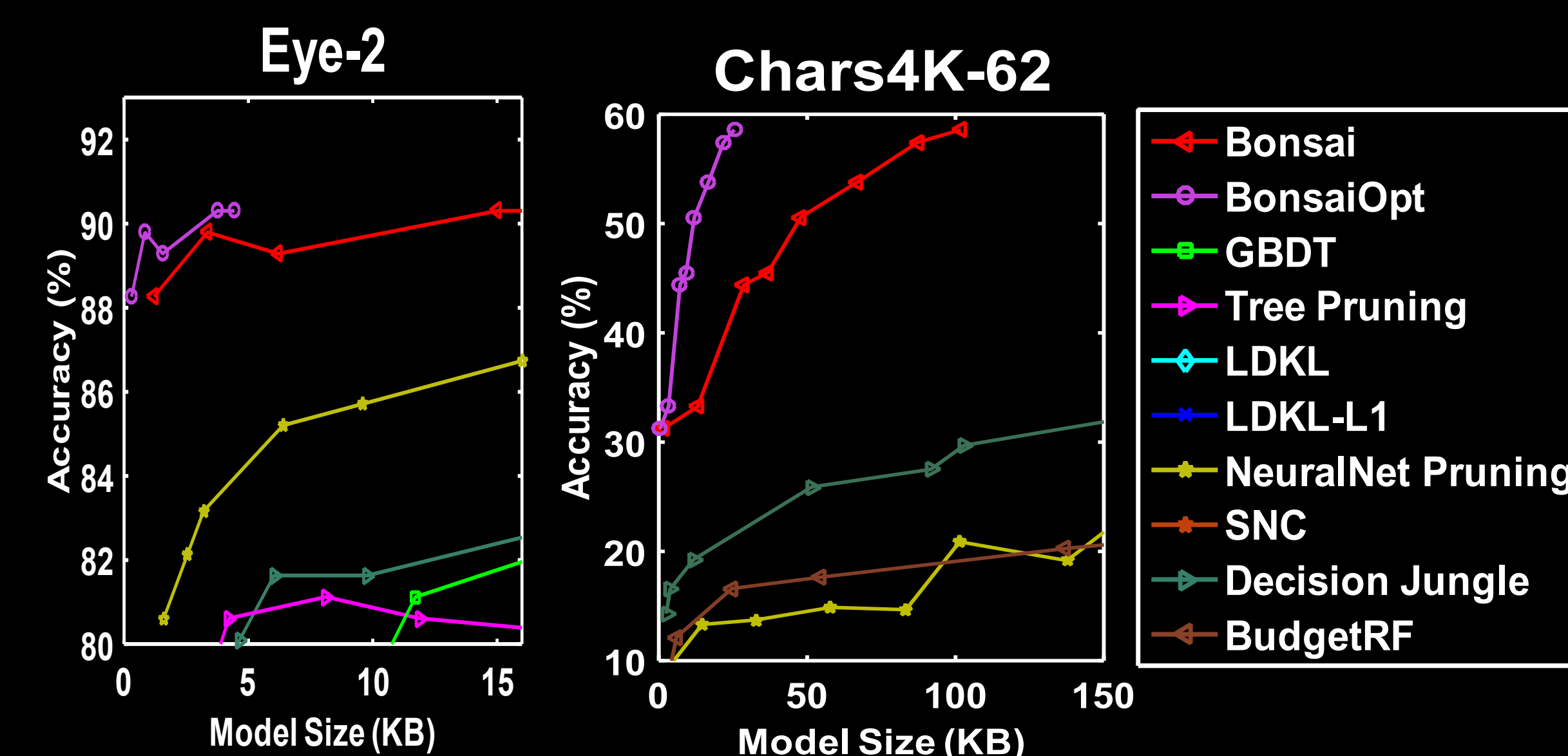$$\Theta_{t+1} = \Theta_t - \eta_\Theta (\nabla_\Theta P)_{supp\{\Theta_t\}}$$
$$Z_{t+1} = Z_t - \eta_Z (\nabla_Z P)_{supp\{Z_t\}}$$

2. Hard thresholding step
$$\Theta_{t+1} = H_B(\Theta_t - \eta_\Theta \nabla_\Theta P)$$
$$Z_{t+1} = H_K(Z_t - \eta_Z \nabla_Z P)$$

## Prediction Accuracy vs Model Size

Eye-2

Chars4K-62

Legend:
- Bonsai
- BonsaiOpt
- GBDT
- Tree Pruning
- LDKL
- LDKL-L1
- NeuralNet Pruning
- SNC
- Decision Jungle
- BudgetRF

## Prediction Accuracy vs Model Size

RTWhale-2    Eye-2    Chars4K-2    WARD-2

MNIST-2    USPS-2    CIFAR10-2

Legend:
- Bonsai
- BonsaiOpt
- GBDT
- Tree Pruning
- LDKL
- LDKL-L1
- NeuralNet Pruning
- SNC
- Decision Jungle
- BudgetRF

## Comparison to Uncompressed Methods

Compressed : Uncompressed :
- Bonsai
- GBDT
- kNN
- RBF-SVM
- Neural Nets

## Prediction Costs on the Arduino Uno

Accuracy (%)    Pred Time (ms) 2187    Pred Energy (mJ) 1312

Eye-2

Eye-2

RTWhale-2    521    313

Legend:
- Bonsai
- NeuralNet
- Linear
- Cloud-GBDT
- LDKL

## Bing L3/L4 Ranker Results

- FastRank
- Bonsai

Code for Bonsai code can be downloaded from http://manikvarma.org/