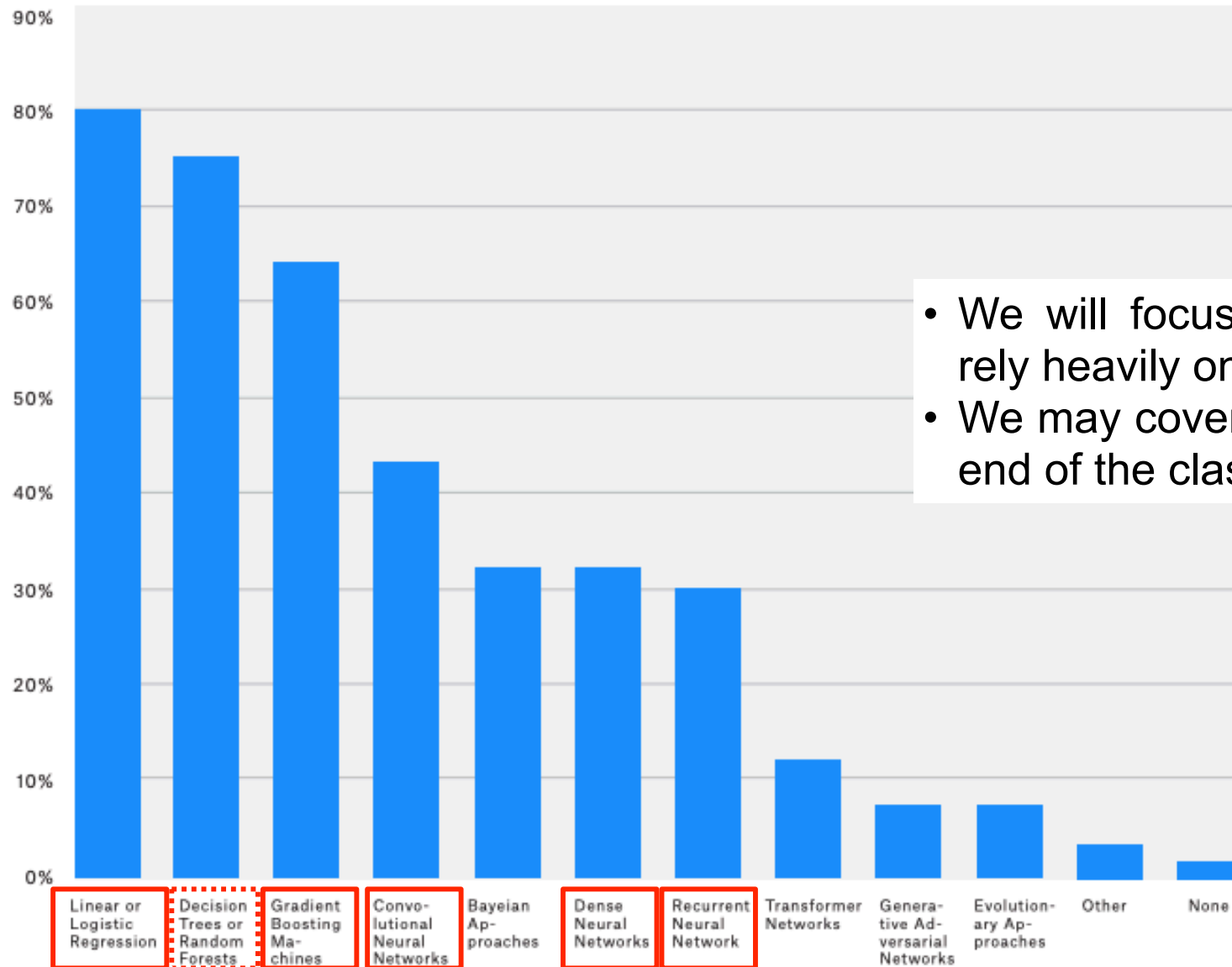


# Decision Trees and Forests

Pascal Fua  
IC-CVLab

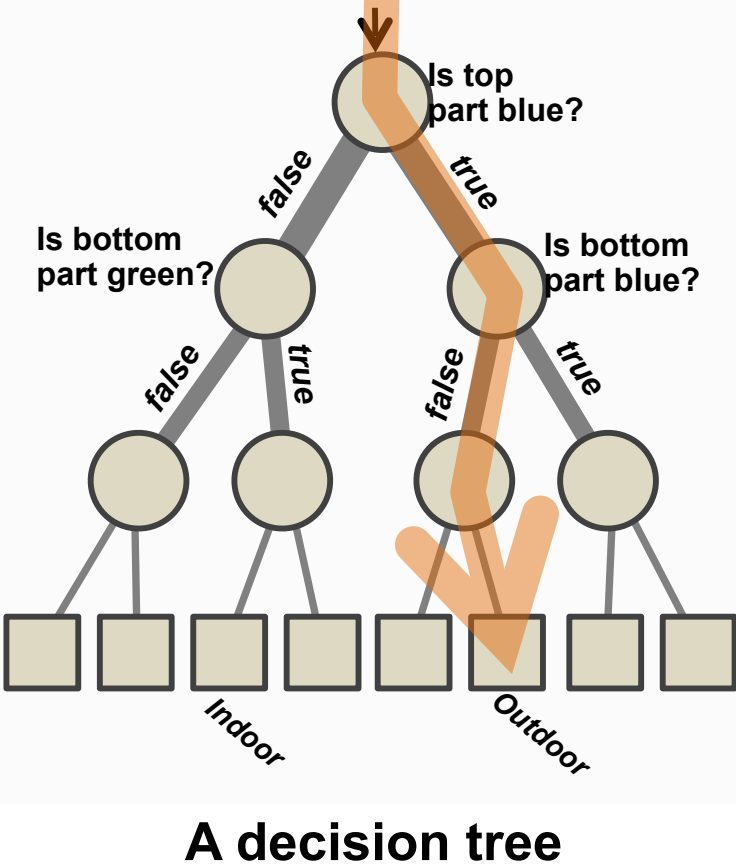
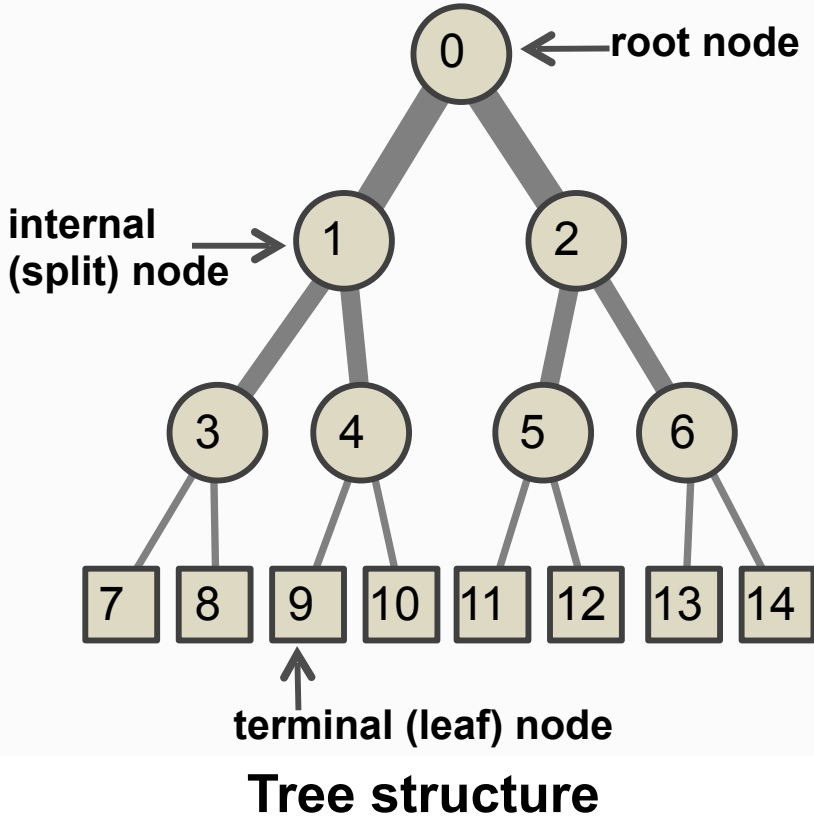
# Kaggle Survey (2019)



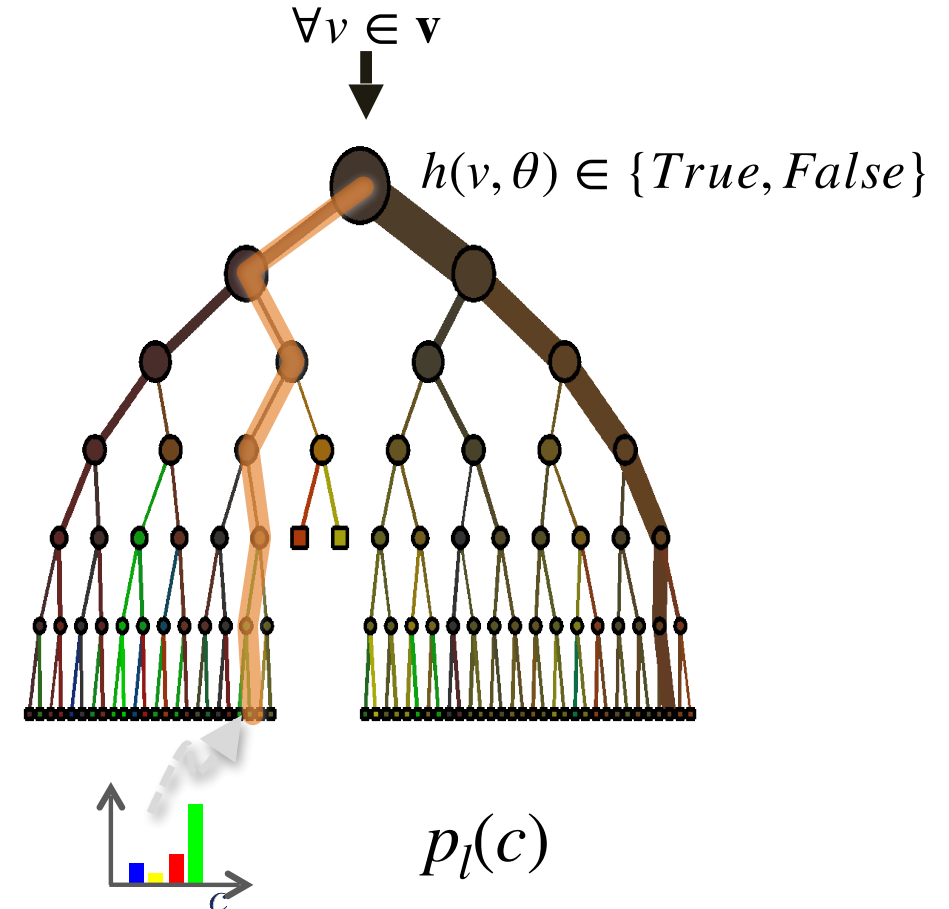
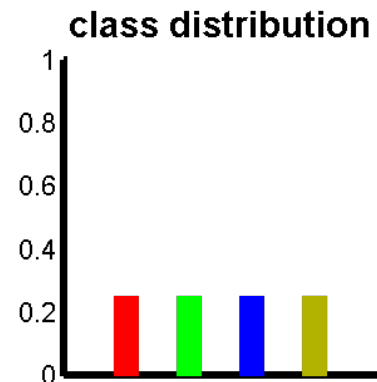
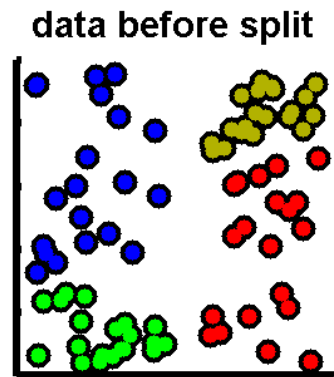
- We will focus on methods that do not rely heavily on probabilities.
- We may cover some of the others at the end of the class time permitting.

What data science methods do you use at work?

# Decision Tree



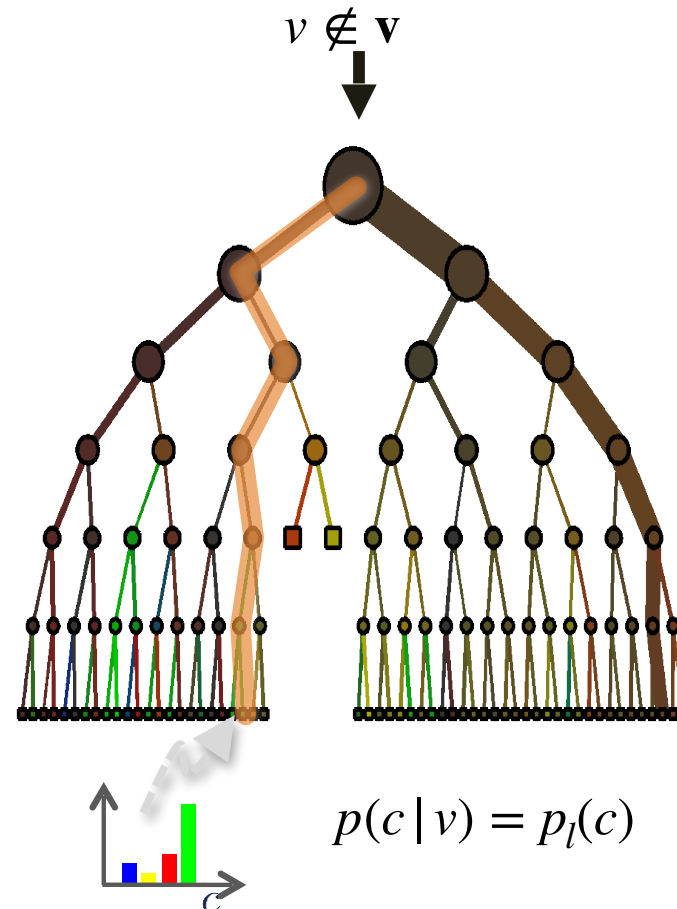
# Training



- Training set  $\mathcal{V}$
- To each sample  $v$  is assigned a class  $c$ .

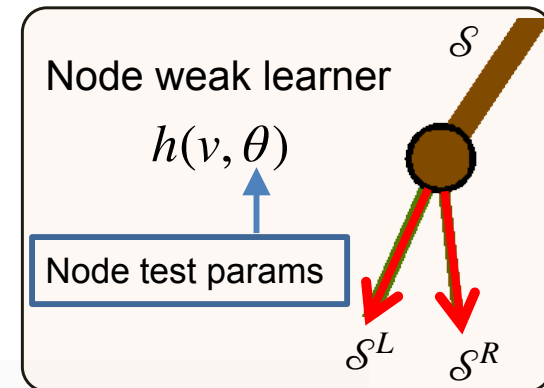
- Compute  $p_l(c)$ , the proportion of samples in each class that lands in leaf 1.

# Testing

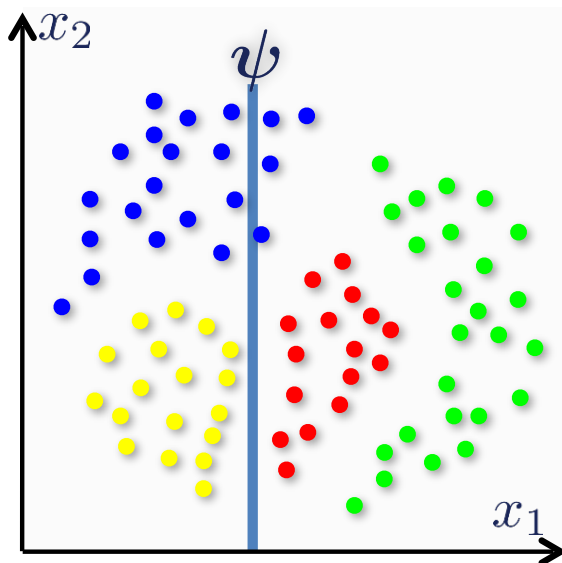


- Let us assume that  $v$  falls into leaf 1.
- We take the probability of belonging to class  $c$ ,  $p(c | v)$ , to be  $p_l(c)$  if it lands in leaf 1.

# Weak Learners



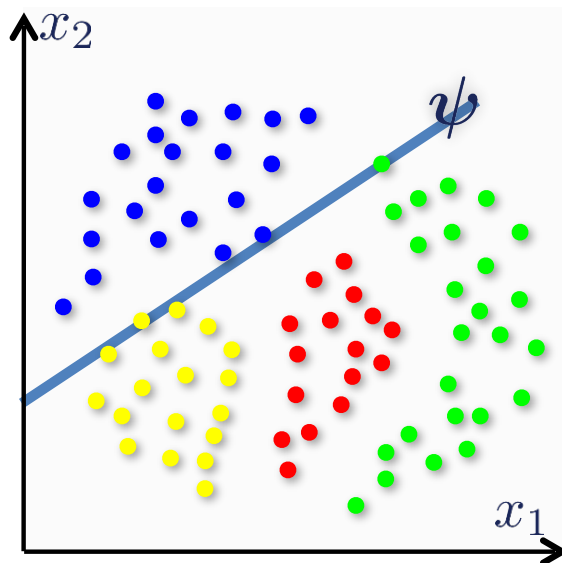
## Weak learner examples



Weak learner: axis aligned.

$$h(\mathbf{v}, \theta) = [\tau_1 > \phi(\mathbf{v}) \cdot \psi > \tau_2]$$

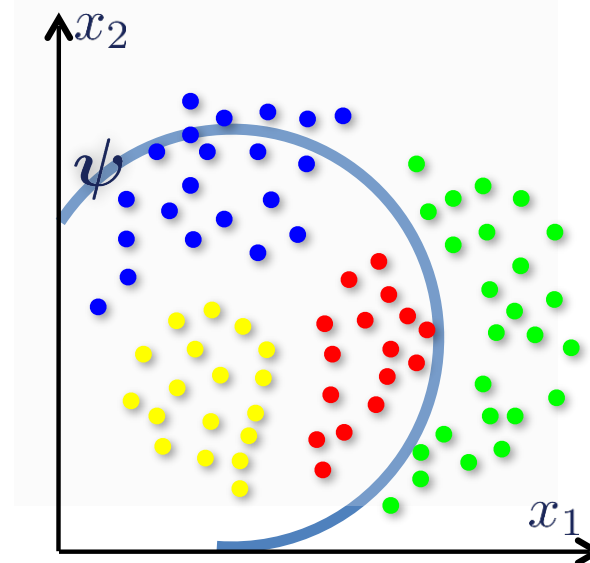
Feature response for 2D example.  $\phi(\mathbf{v}) = (x_1 \ x_2 \ 1)^\top$   
 With  $\psi = (1 \ 0 \ \psi_3)$  or  $\psi = (0 \ 1 \ \psi_3)$



Weak learner: oriented line.

$$h(\mathbf{v}, \theta) = [\tau_1 > \phi(\mathbf{v}) \cdot \psi > \tau_2]$$

Feature response for 2D example.  $\phi(\mathbf{v}) = (x_1 \ x_2 \ 1)^\top$   
 With  $\psi \in \mathbb{R}^3$  a generic line in homog. coordinates.



Weak learner: conic section.

$$h(\mathbf{v}, \theta) = [\tau_1 > \phi^\top(\mathbf{v}) \psi \phi(\mathbf{v}) > \tau_2]$$

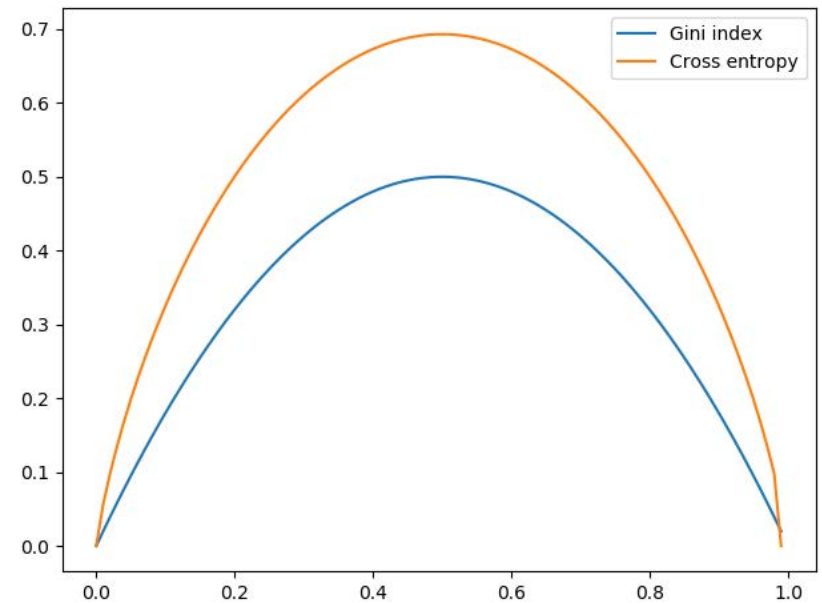
Feature response for 2D example.  $\phi(\mathbf{v}) = (x_1 \ x_2 \ 1)^\top$   
 With  $\psi \in \mathbb{R}^{3 \times 3}$  a matrix representing a conic.

# Entropy and Gini Index

Let  $p^k$  be the proportion of data points in  $\mathcal{S}$  that are assigned to class  $k$ .  
We can define

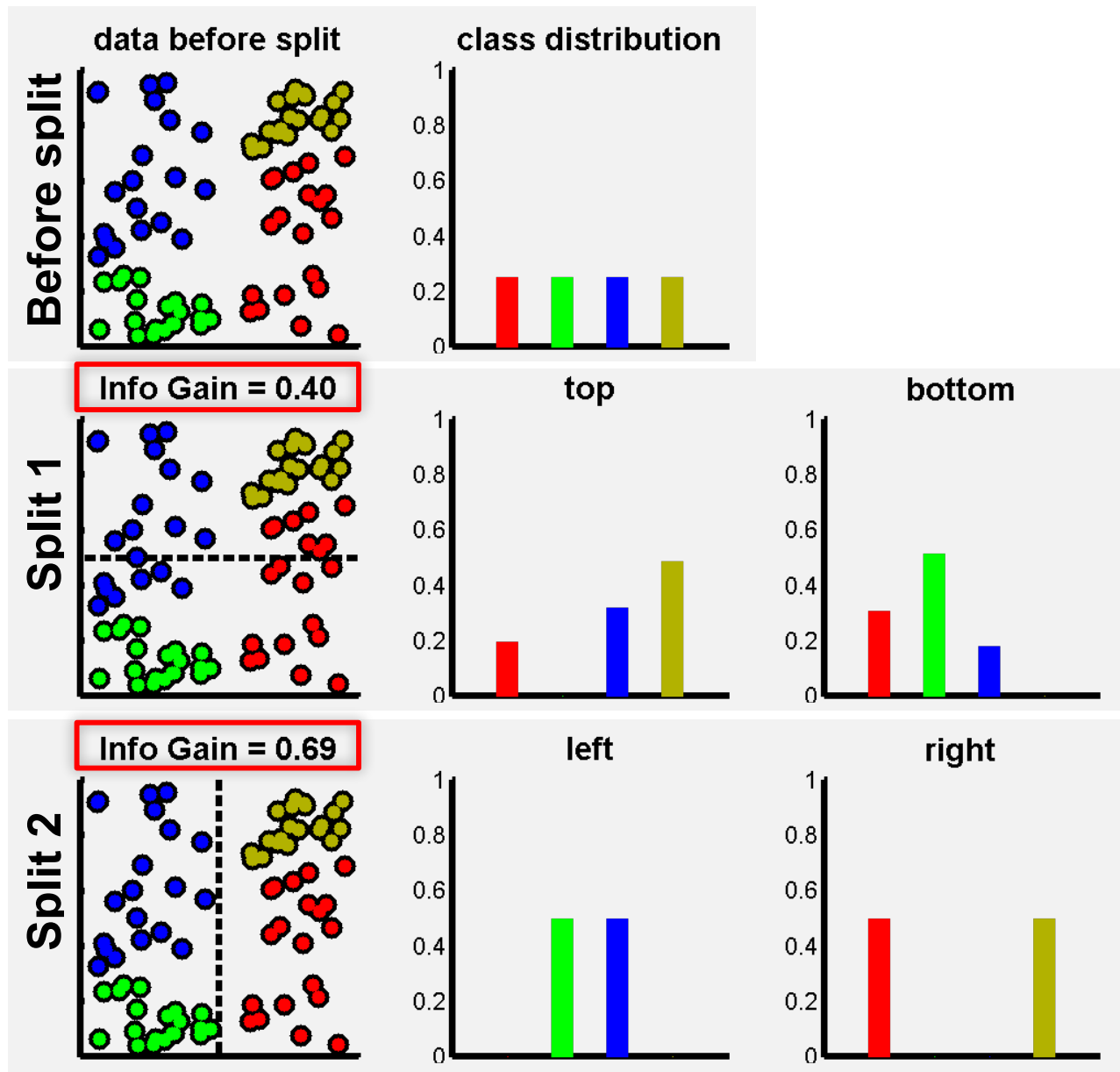
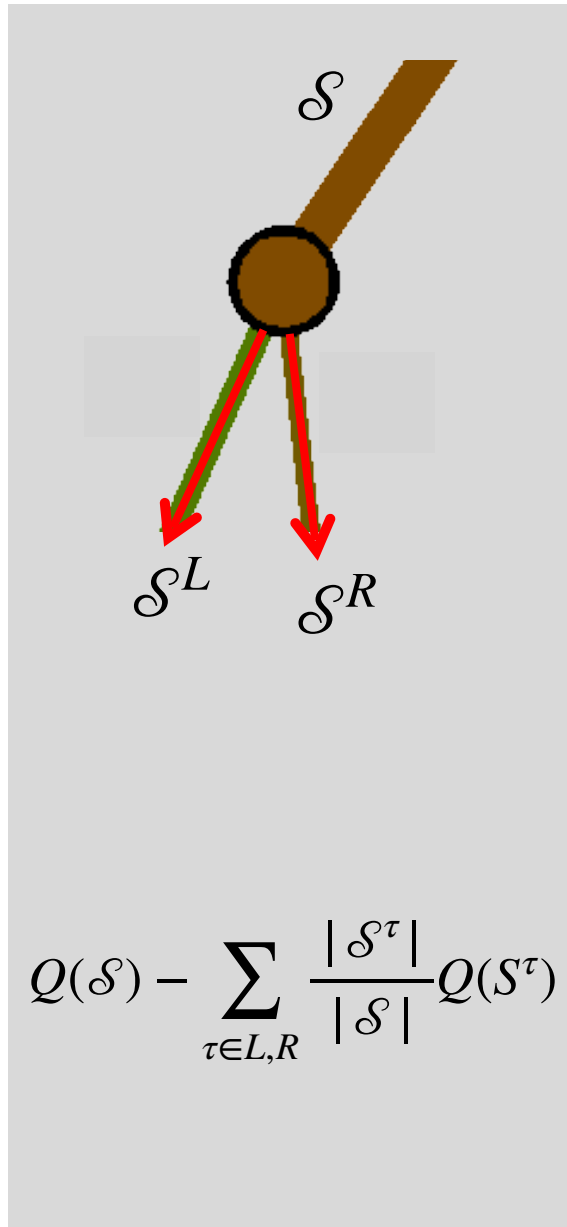
- the Gini index  $Q(\mathcal{S}) = \sum_{k=1}^K p^k(1 - p^k)$ ,
- the entropy  $Q(\mathcal{S}) = - \sum_{k=1}^K p^k \ln p^k$ .
- They both vanish when  $\exists k, p^k = 1$ .
- They are maximized when all  $p^k$  are equal.

—> Minimizing these measures favors leaves in which a large fraction of samples belong to the same class.



Two classes case.

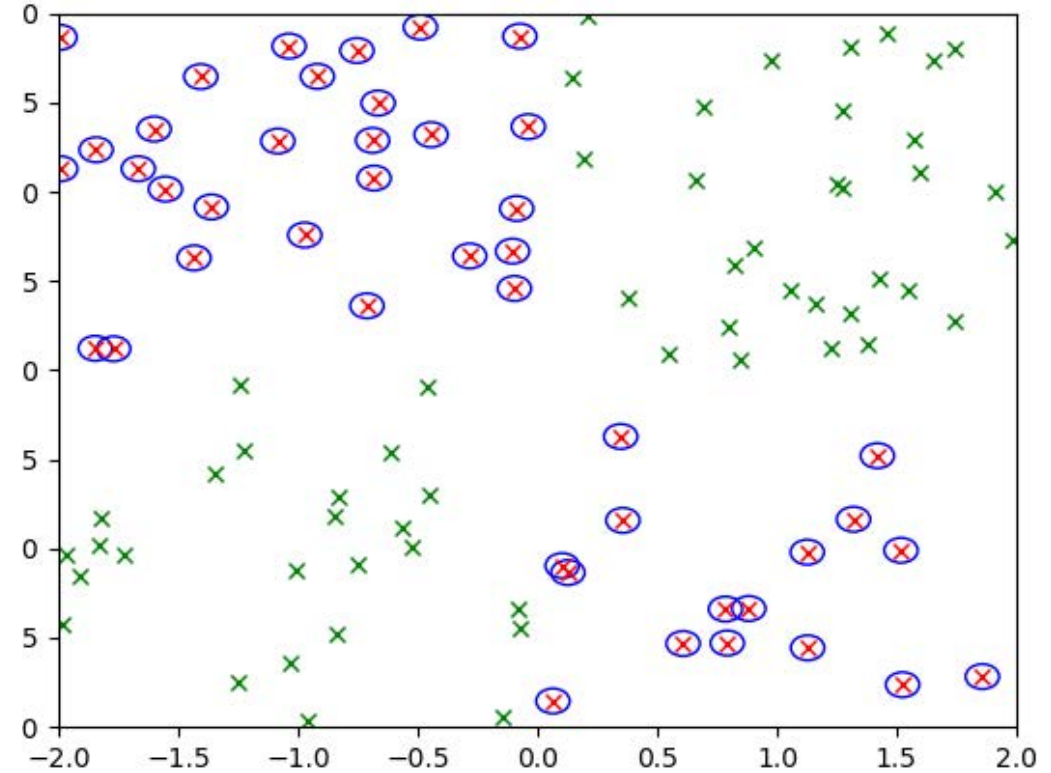
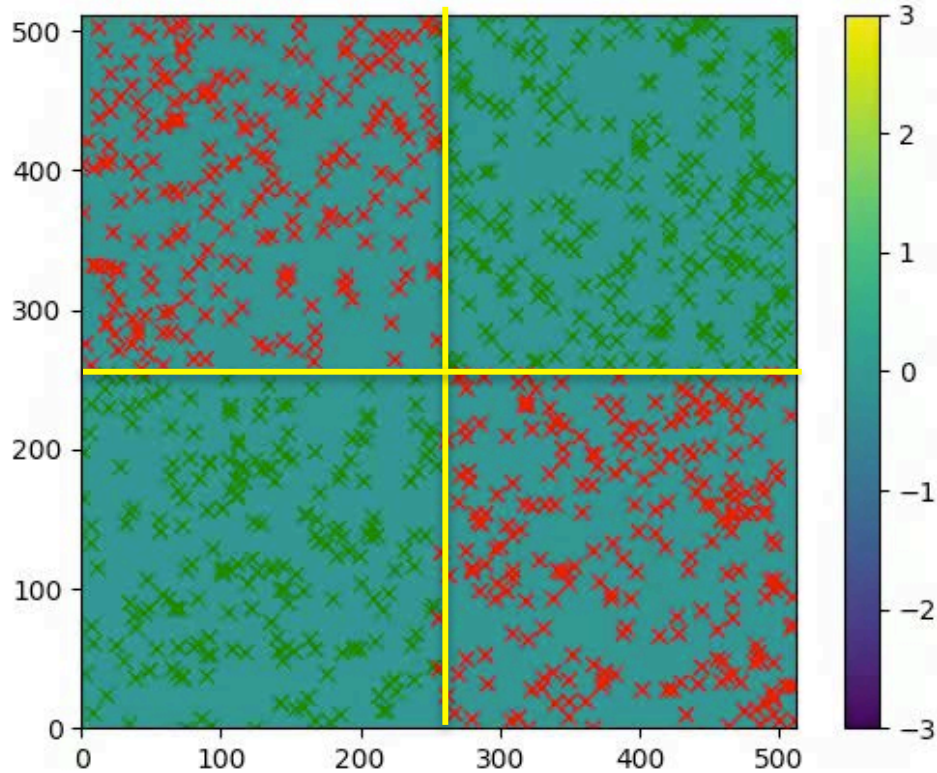
# Maximizing Information Gain



At each node, pick the weak learner that delivers the highest information gain.



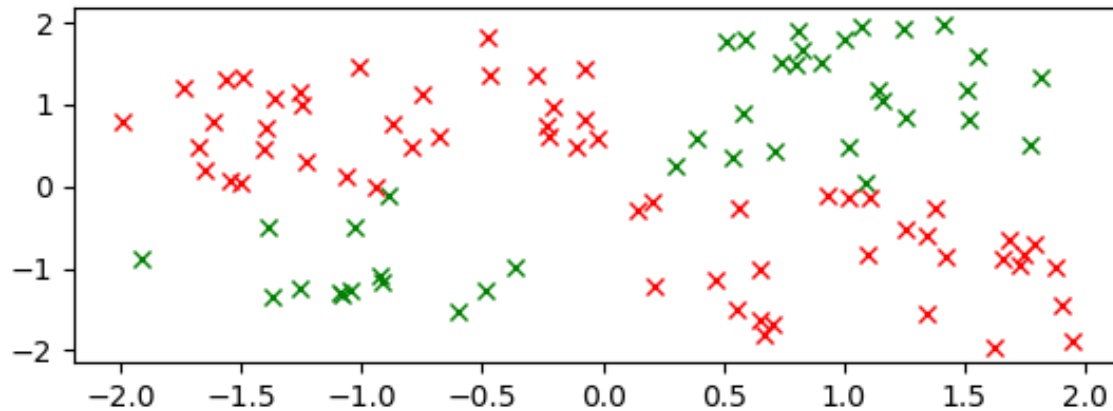
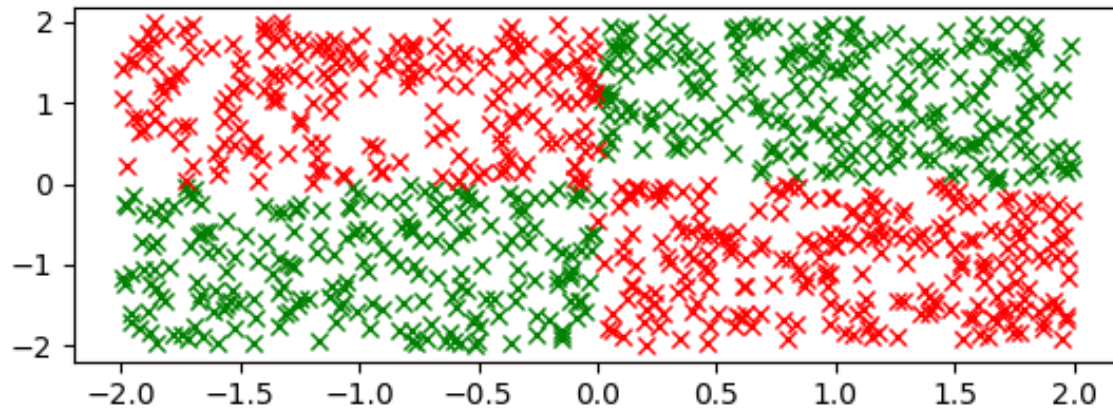
# Problematic for AdaBoost ....



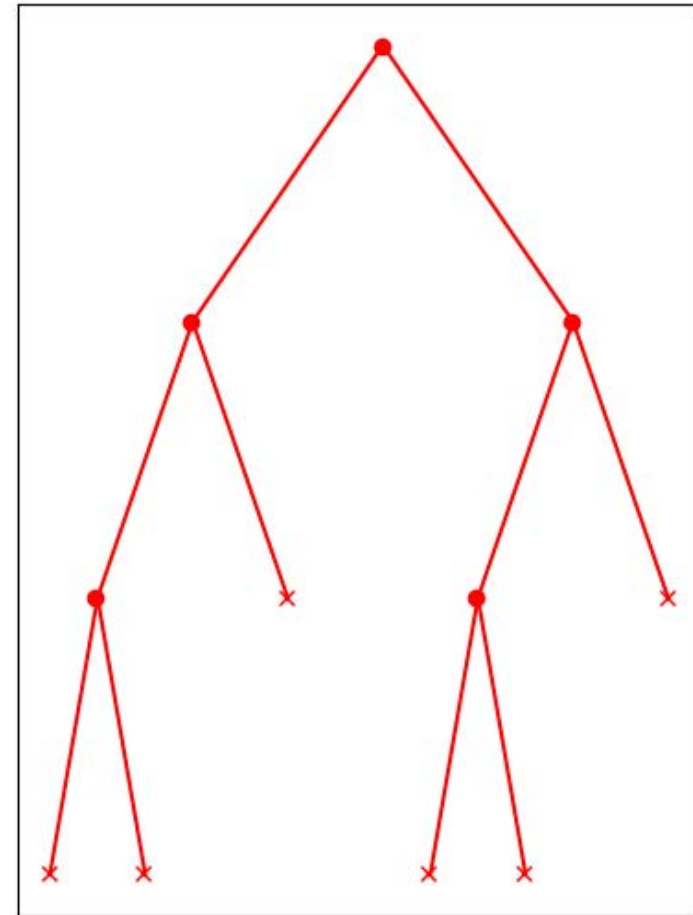
When using linear classifiers as weak learners.

# ... but not for Trees

Training

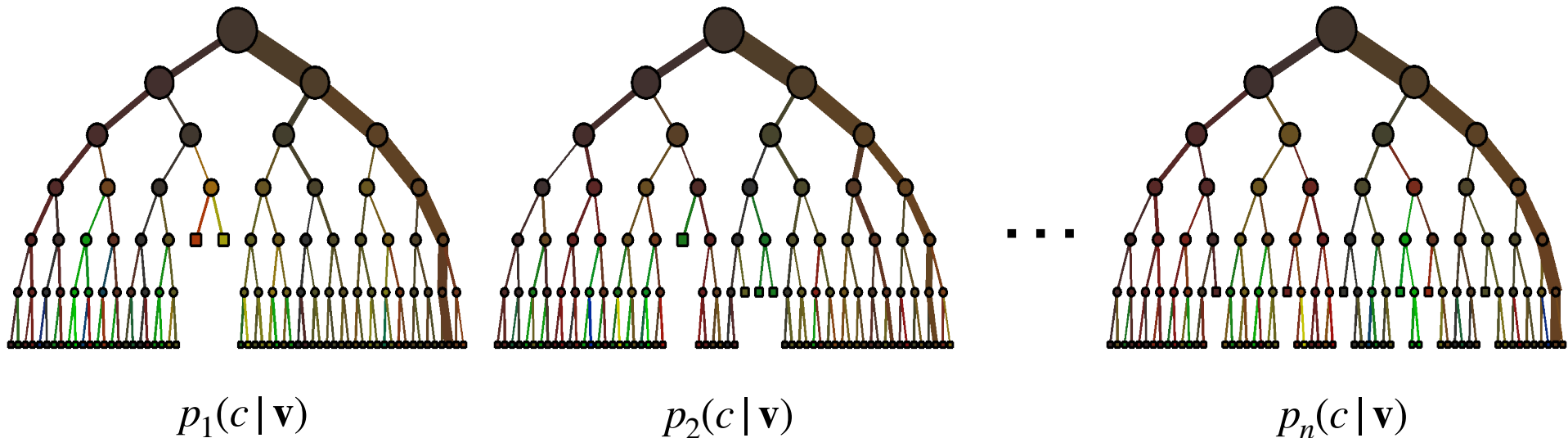


Validation



# From Trees to Forests

Use multiple trees to increase robustness:



$$p(c|\mathbf{v}) = f(p_1(c|\mathbf{v}), \dots, p_T(c|\mathbf{v}))$$

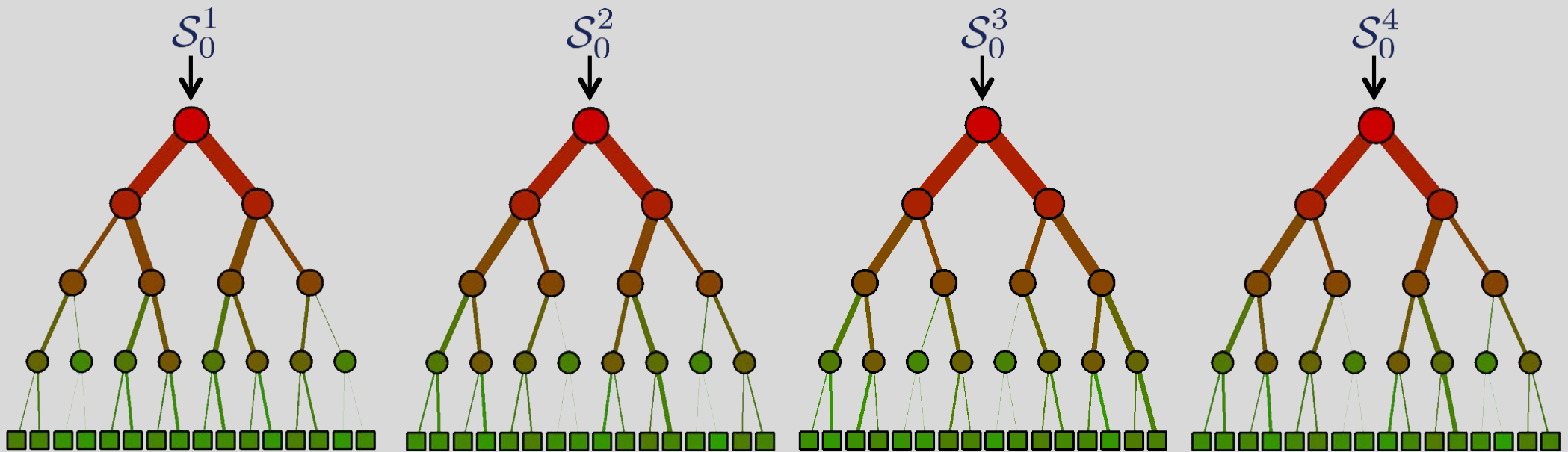
- How many trees?
- How different should they be?
- How do we fuse their outputs?

# Creating Multiple Trees

$\mathcal{S}_0$  Full training set

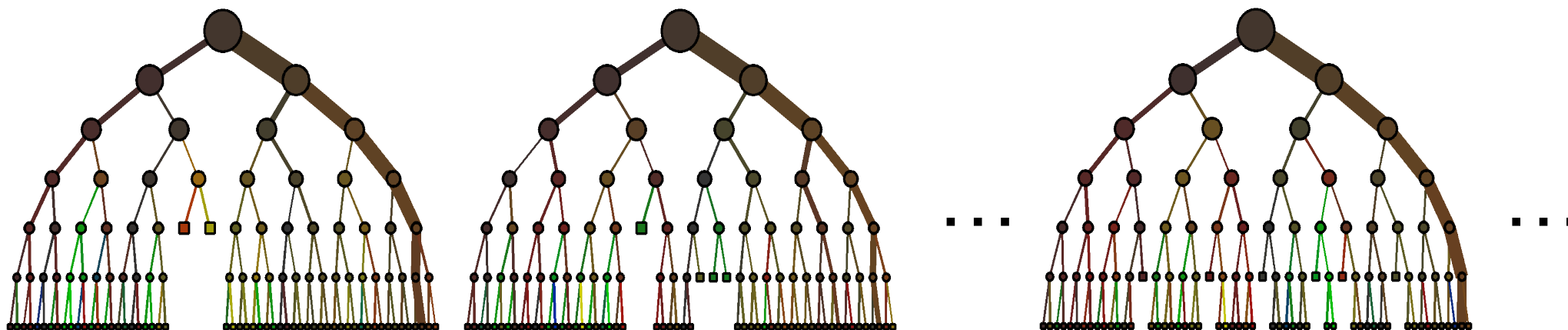
$\mathcal{S}_0^t \subset \mathcal{S}_0$  Randomly sampled subsets made available to train the tree  $t$

Forest training



- The subsets are typically chosen randomly with replacement.
- This is known as bagging.

# Fusing the Output



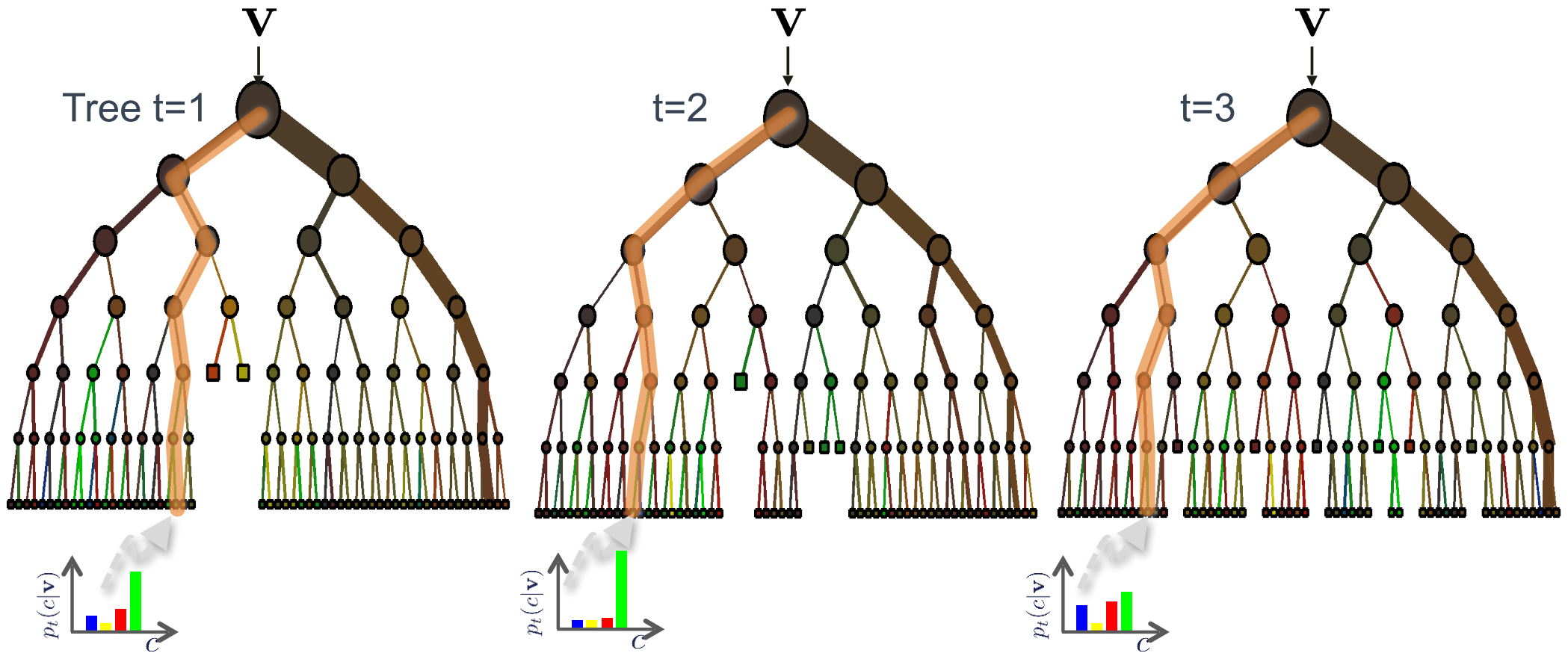
Naive Bayesian:

$$p(c | \mathbf{v}) \propto \prod_t p_t(c | \mathbf{v})$$

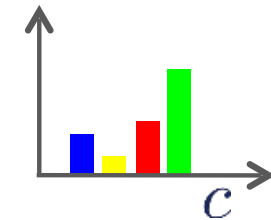
$$L(c, \mathbf{v}) = \frac{1}{T} \sum_t -\log(p_t(c | \mathbf{v}))$$

- Assumes the output of each tree is independent from each other.
- Valid assumption if the training subsets are disjoint.
- Justifiable assumption if the training database is large enough.

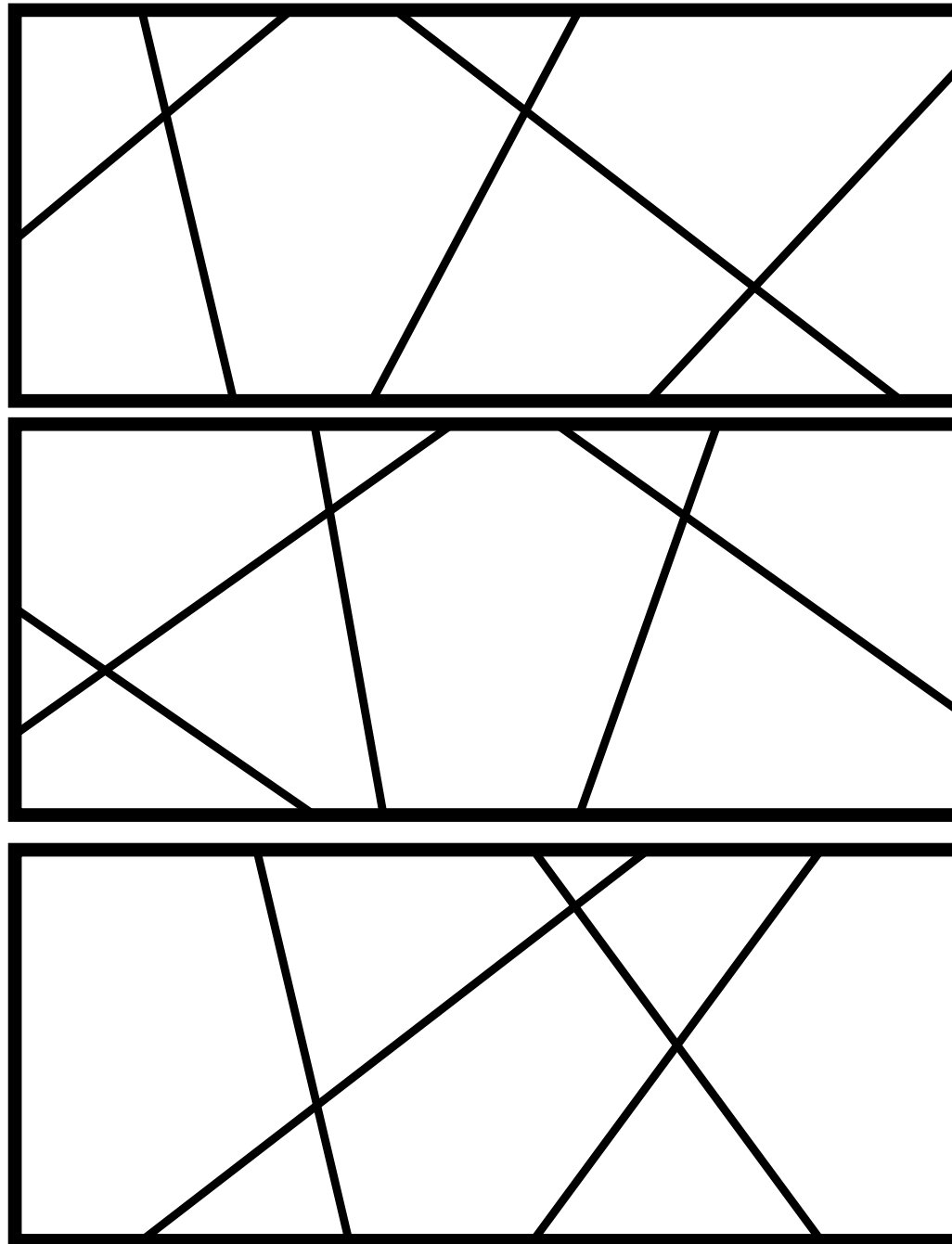
# Ensemble Model



$$L(c, \mathbf{v}) = \frac{1}{T} \sum_t -\log(p_t(c|\mathbf{v}))$$

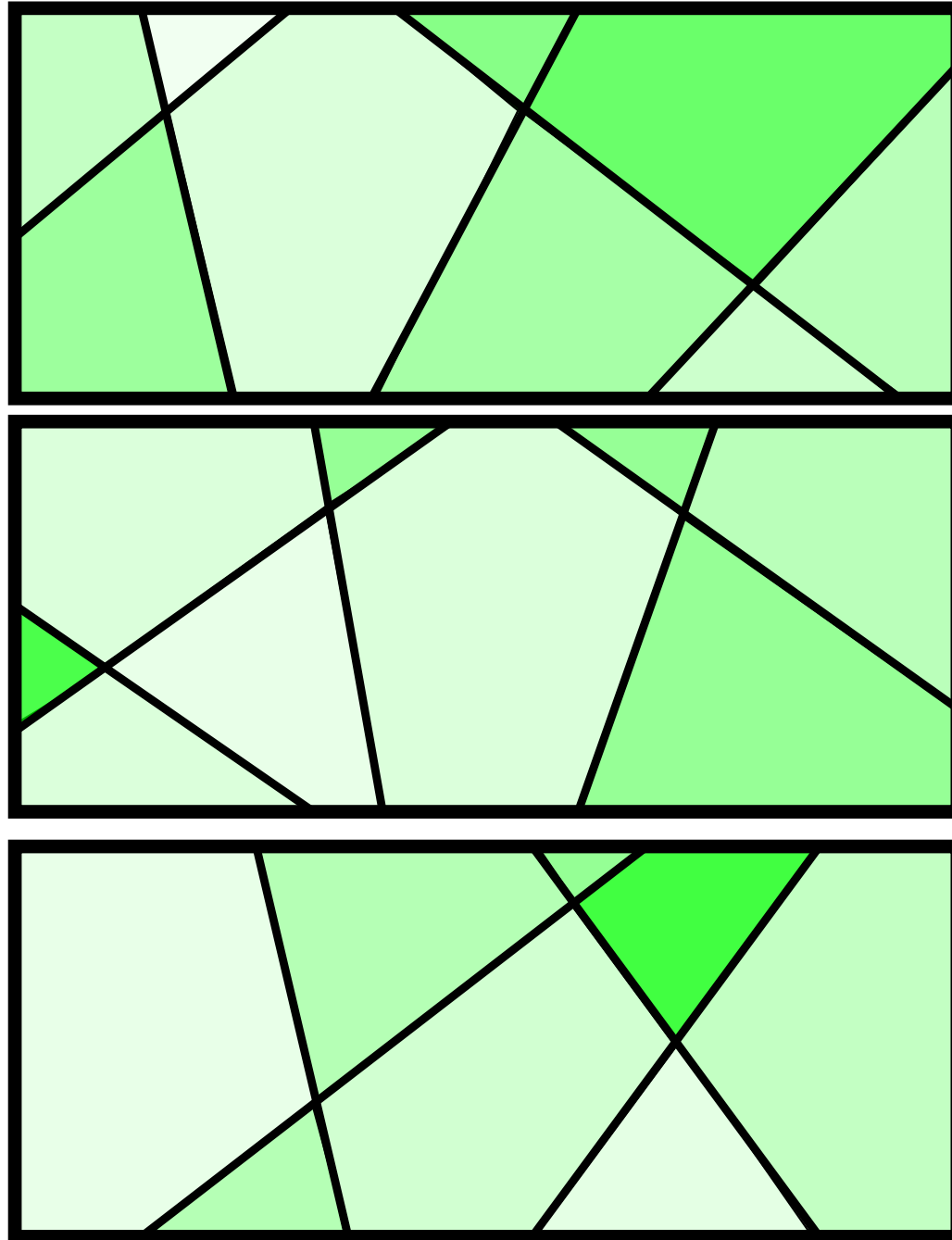


# Graphical Interpretation



Weak classifiers  
at every level of  
the tree split  
the space.

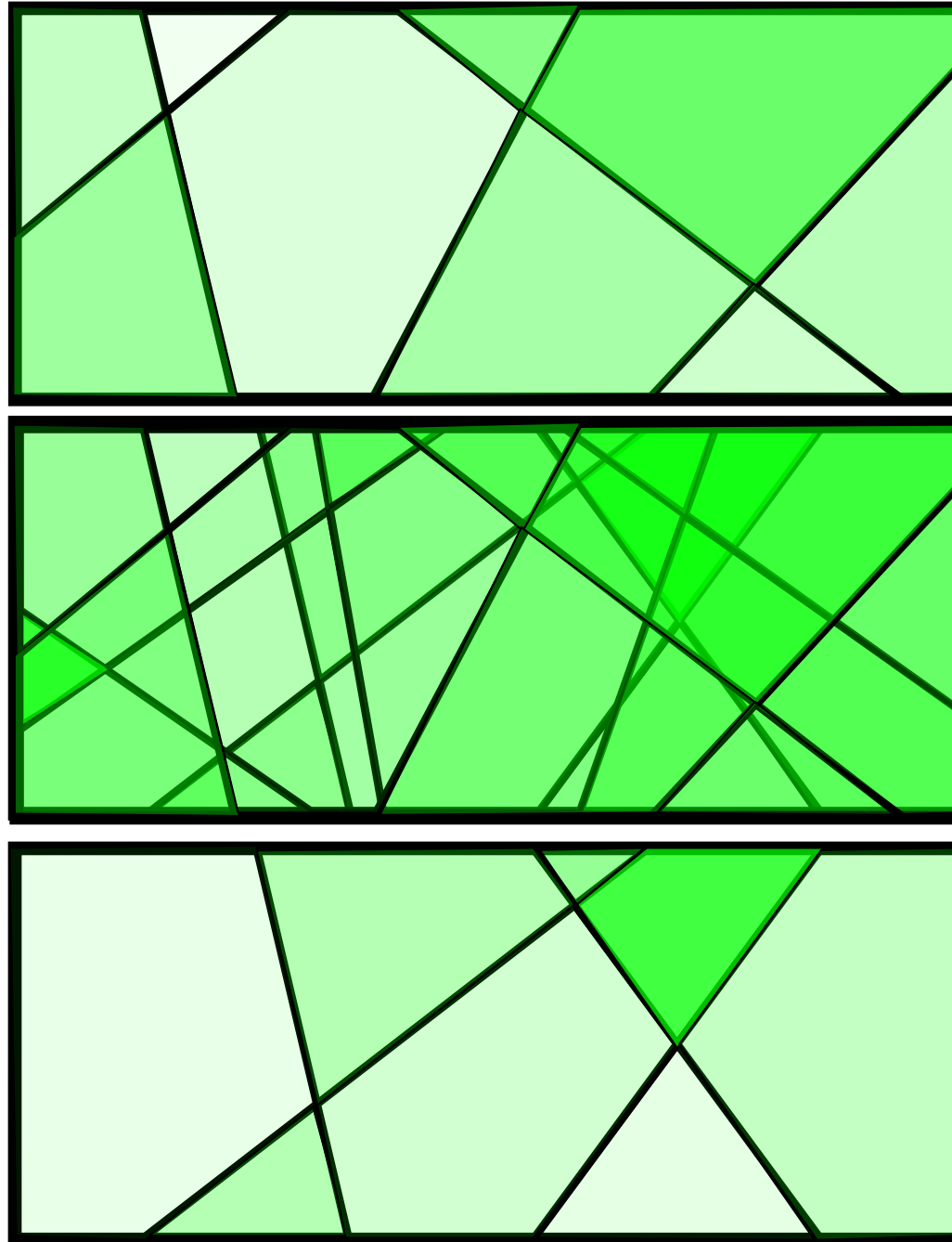
# Graphical Interpretation



Weak classifiers at every level of the tree split the space.

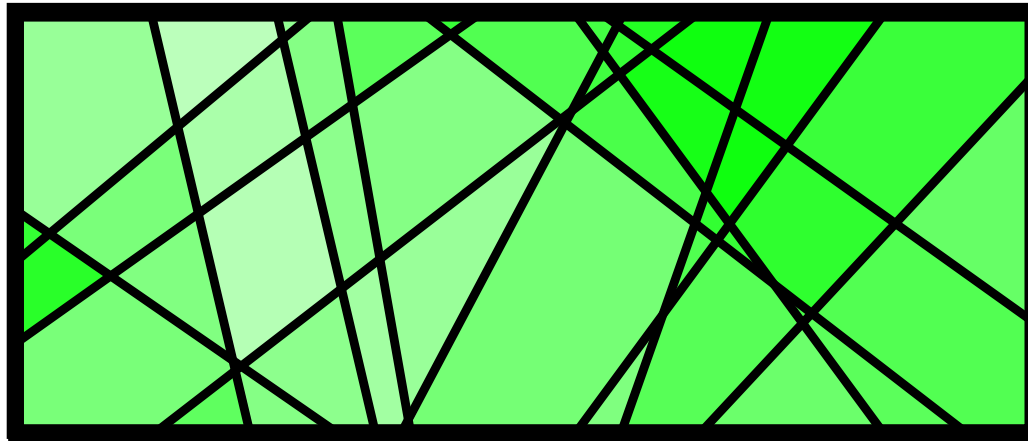


# Graphical Interpretation



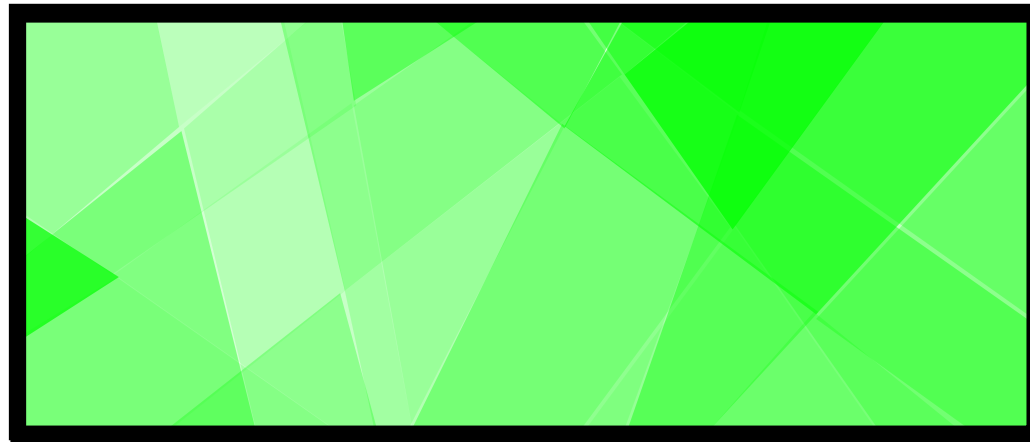
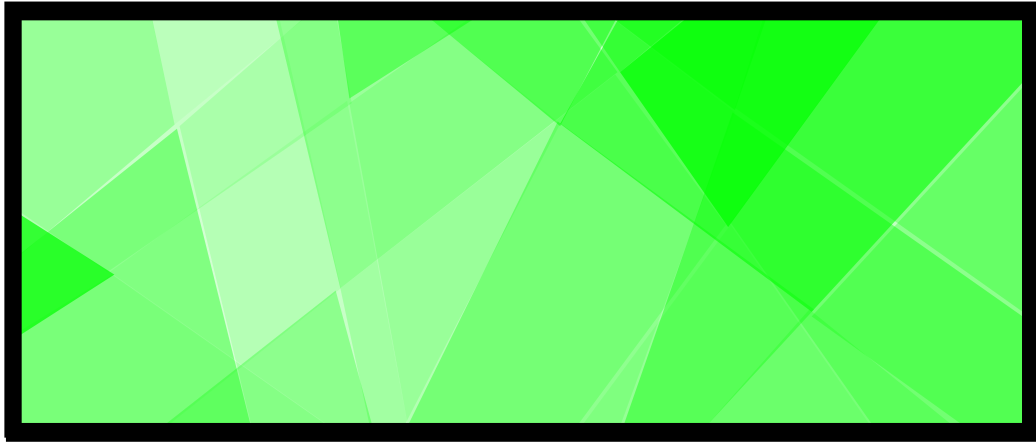
The splits are combined by the hierarchical nature of the tree.

# Graphical Interpretation

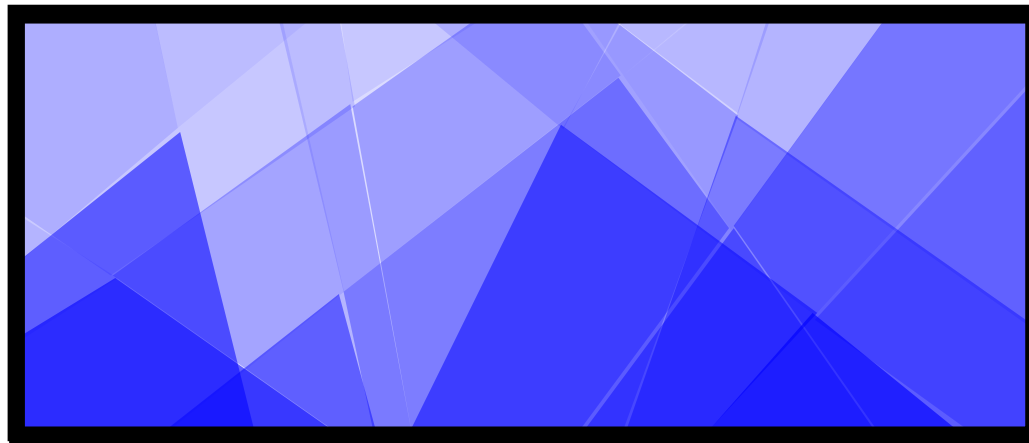
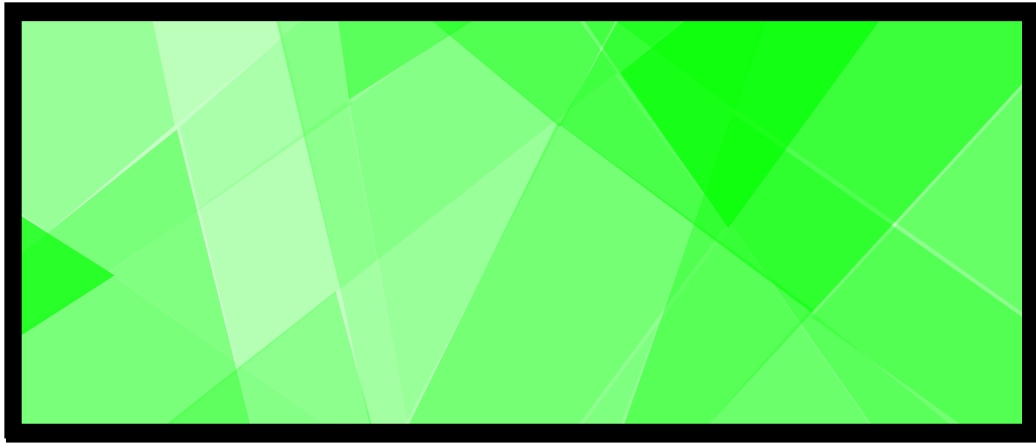


The splits are combined by the hierarchical nature of the tree.

# Graphical Interpretation



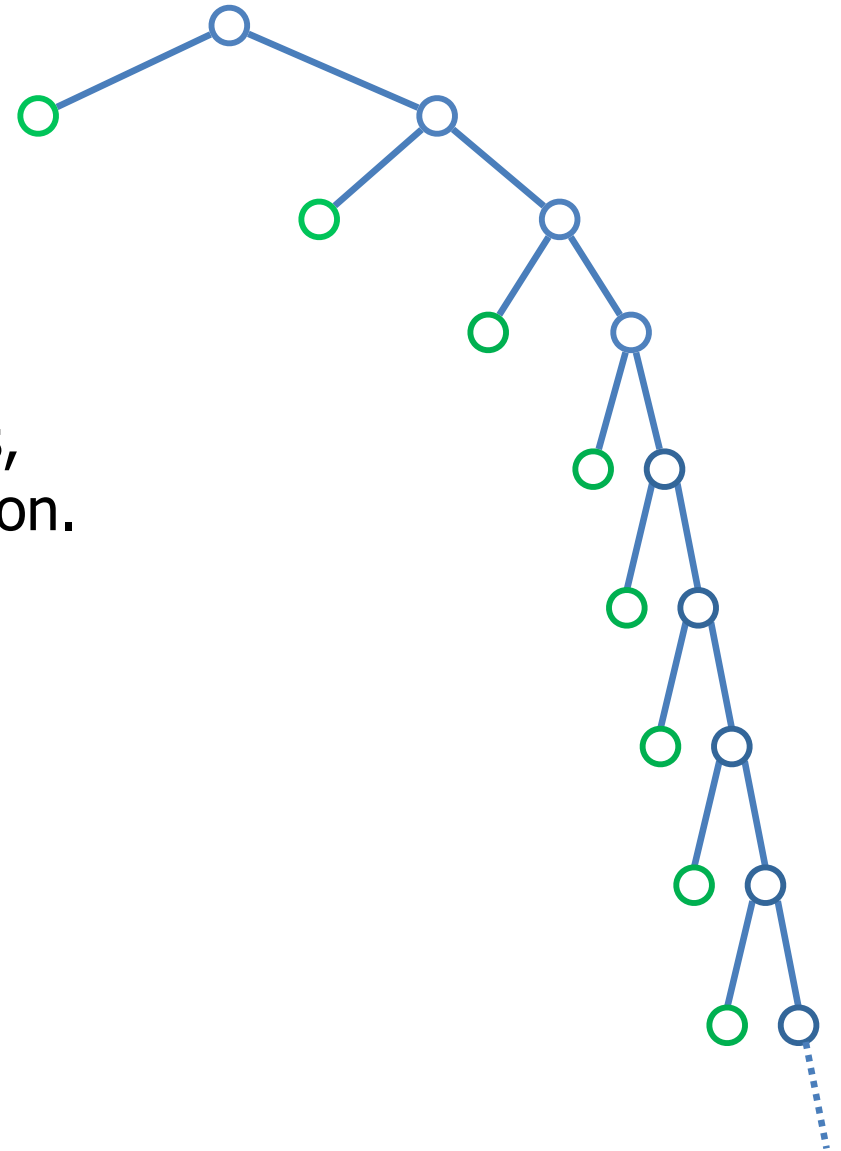
# Graphical Interpretation



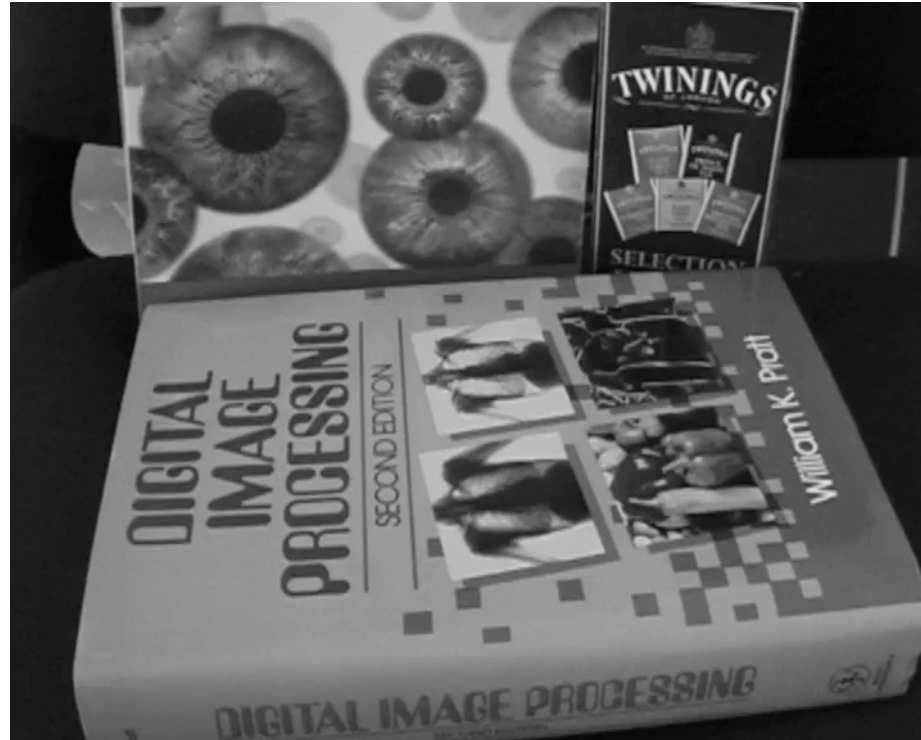
- Each tree produces its own partition of the space.
- These partitions are combined in a Naive Bayesian manner.

# Relationship to Boosting

- **Boosted Cascades:**
  - Very unbalanced tree.
  - Good for unbalanced binary problems, such as sliding window object detection.
- **Randomized forests:**
  - Less deep, more balanced.
  - Ensemble of trees gives robustness.
  - Good for multi-class problems.



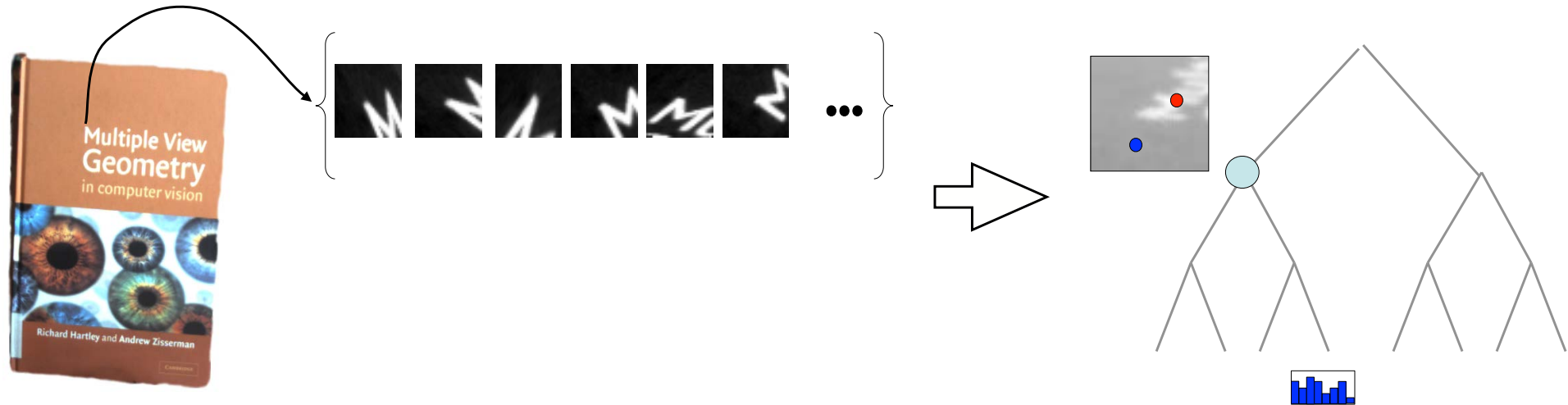
# 3D Pose Estimation



To track the car:

1. We track interest points in the image.
2. We infer their 3D position from the tracks.

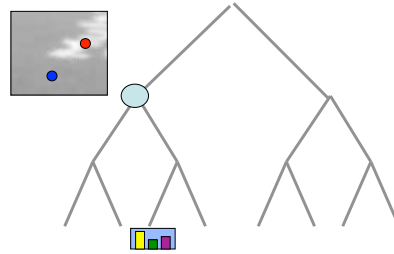
# Classification-Based Approach to Matching



- One class per keypoint.
- Train a decision forest to recognize them.

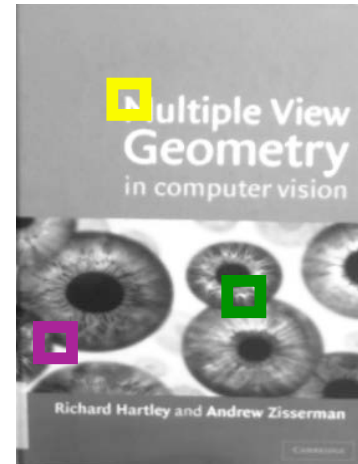
# Simple Weak Learners

The nodes contain simple tests of the form “Is  $I(\mathbf{m}_1) > I(\mathbf{m}_2)$  ?”



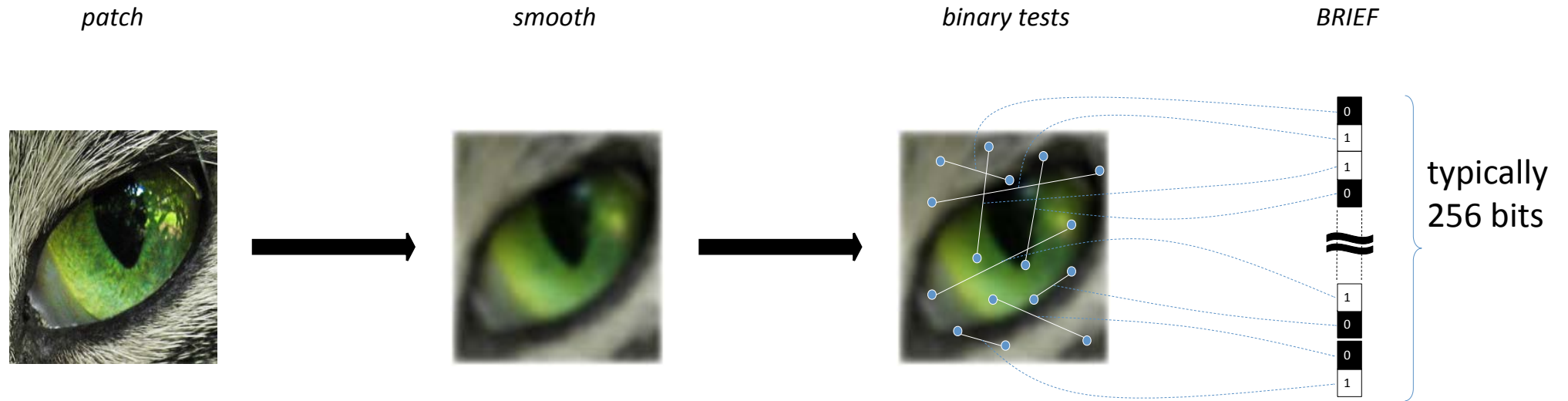
Posteriors can be learned from:

- Warped images
- Video sequences



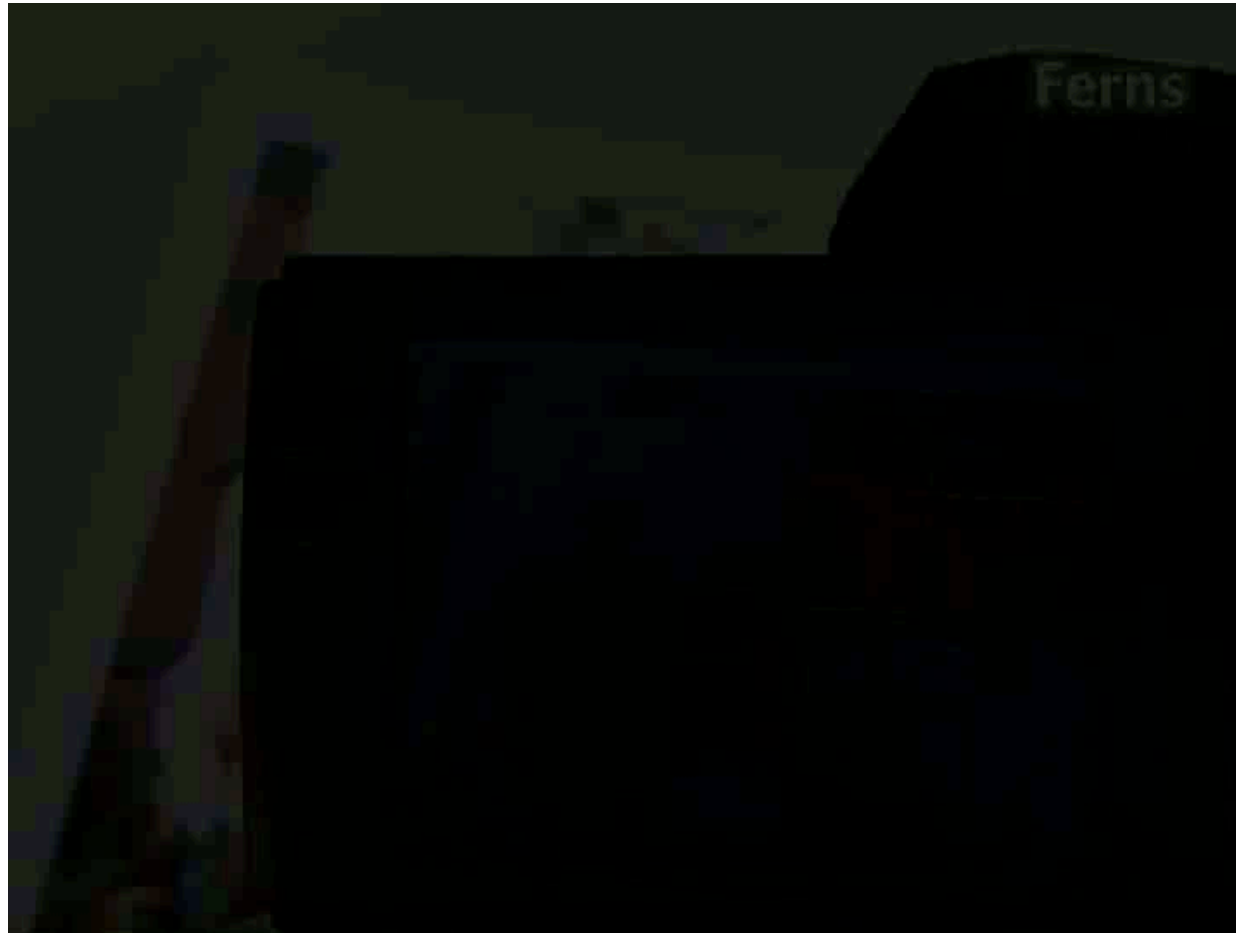


# BRIEF



- Most smooth kernels work, even simple box filters.
- 128, 256, or 512 binary tests usually suffice.
- Random arrangement of tests effective as long as they are evenly sampled.

# Point Correspondences



--> Real-time on a 2008 cell phone.

# Body Part Estimation

depth sensor

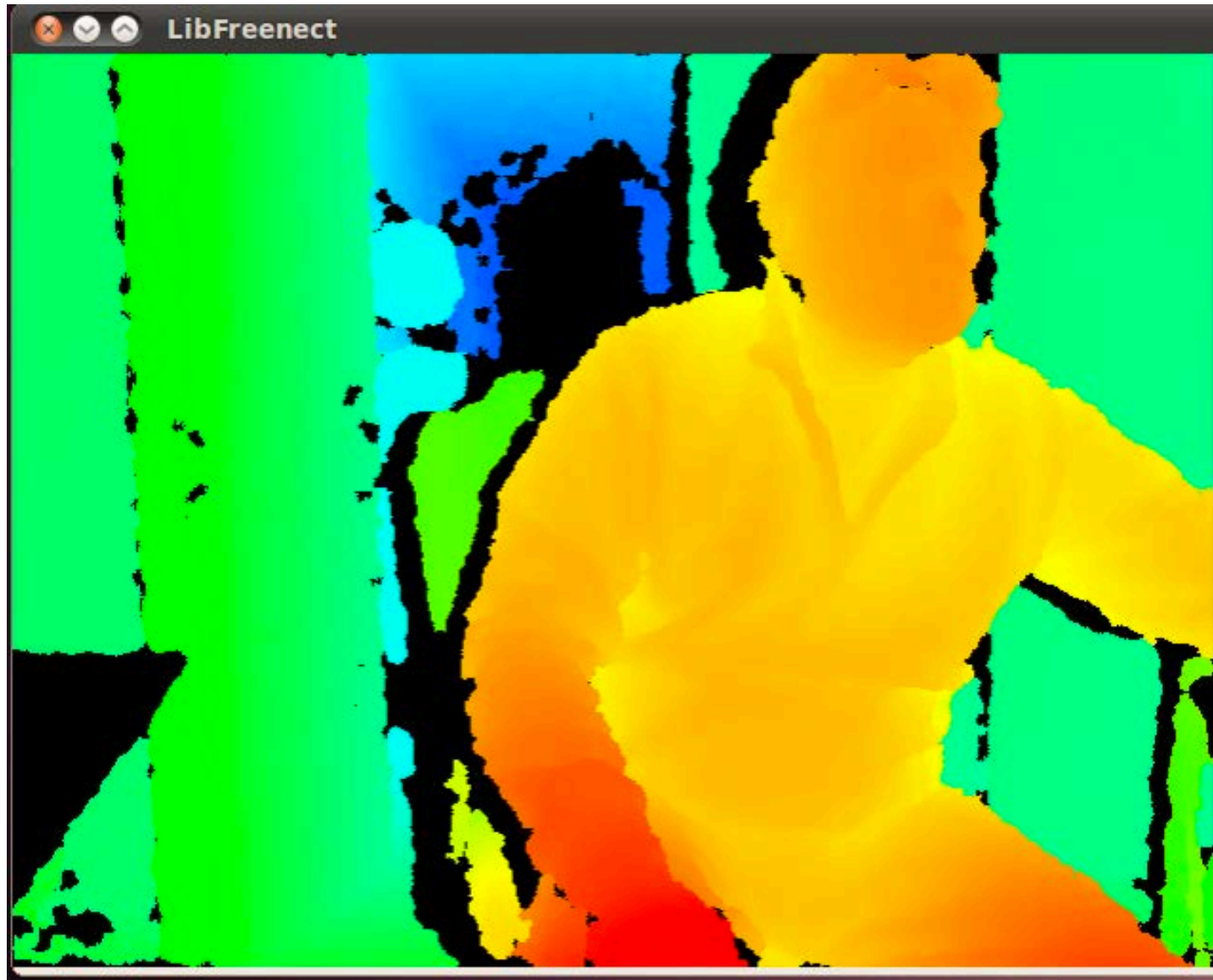
infrared  
emitter

infrared  
camera

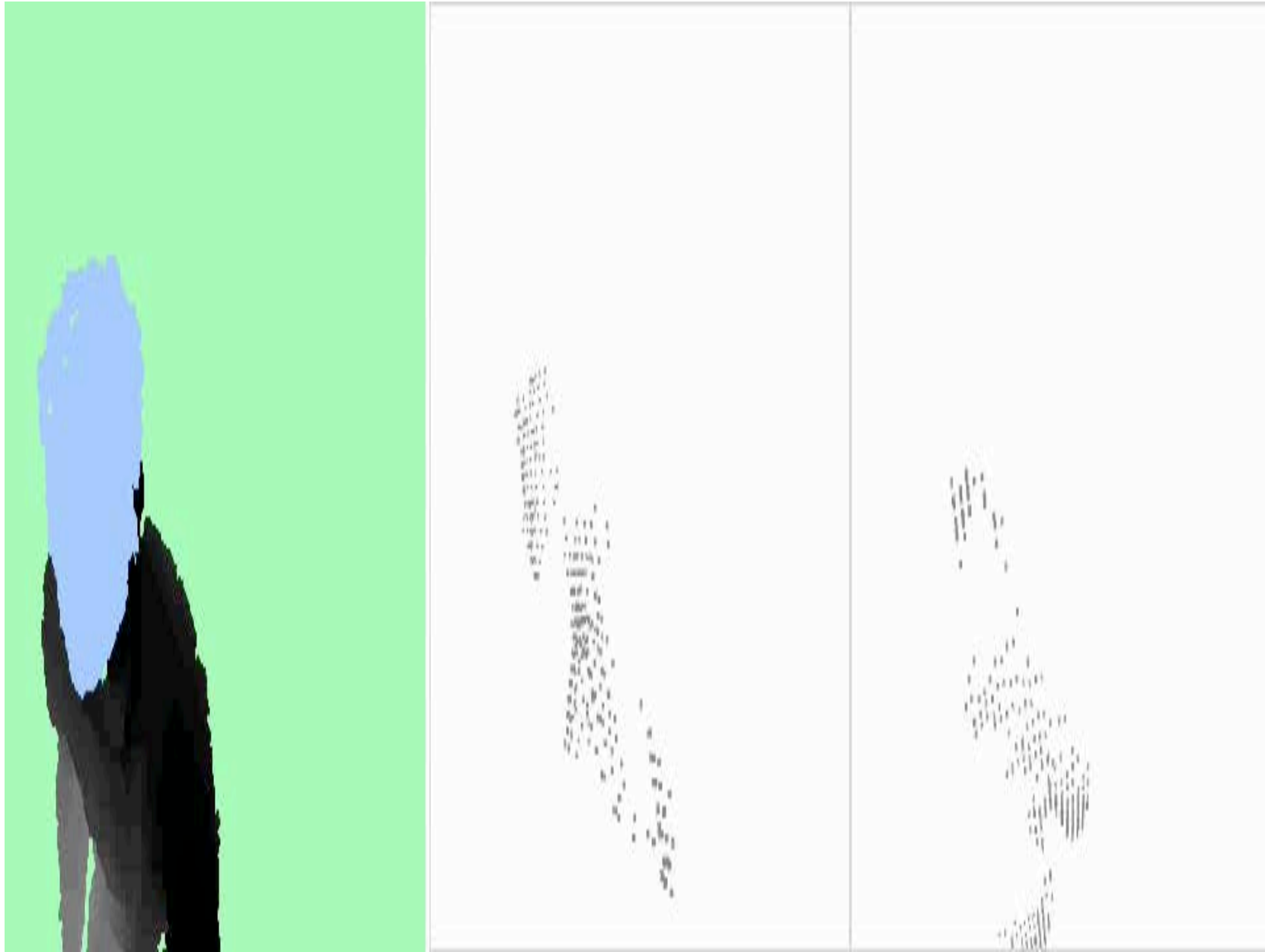
RGB  
camera



# Depth Image



# Depth Sequence



Depth image.

Side view

Top view

# Processing Pipeline

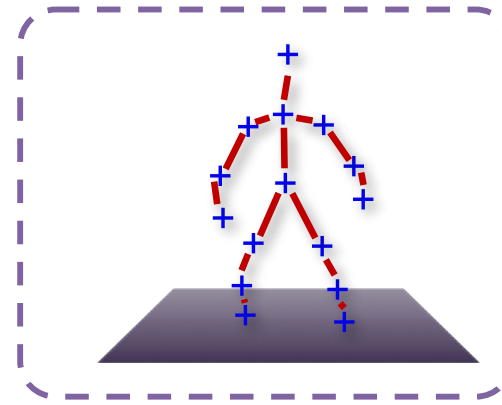
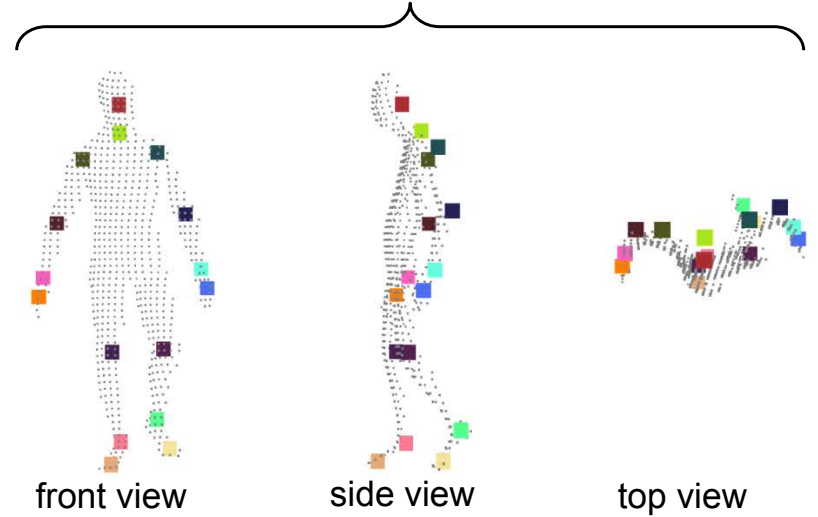
input depth image



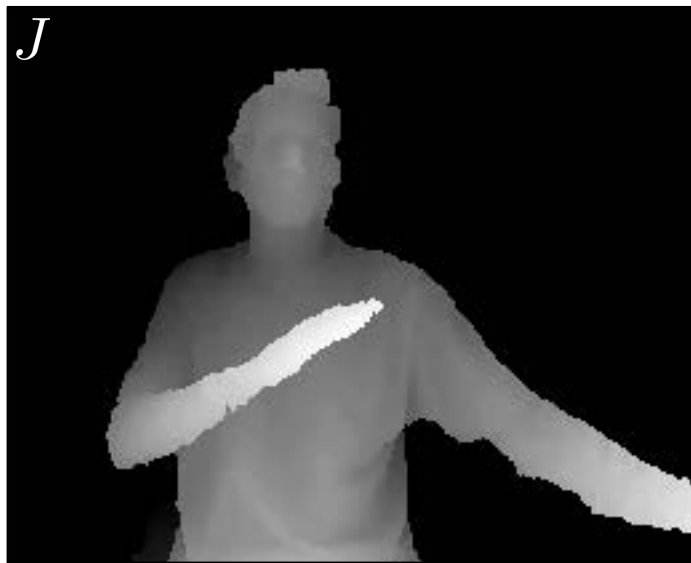
body parts



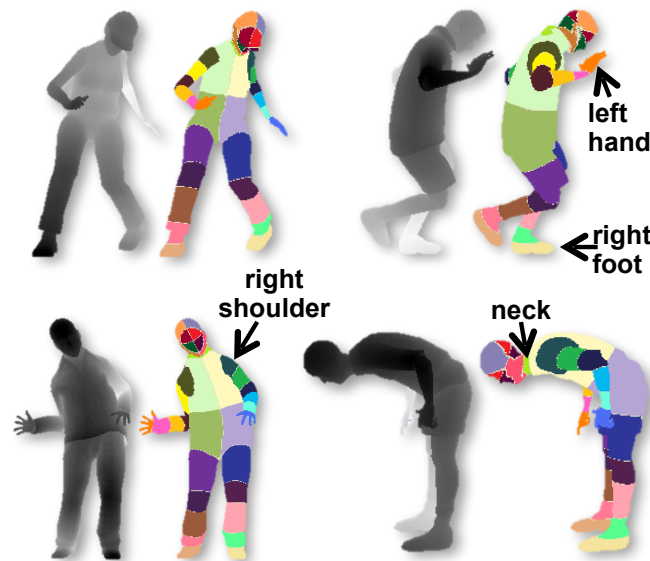
body joint hypotheses



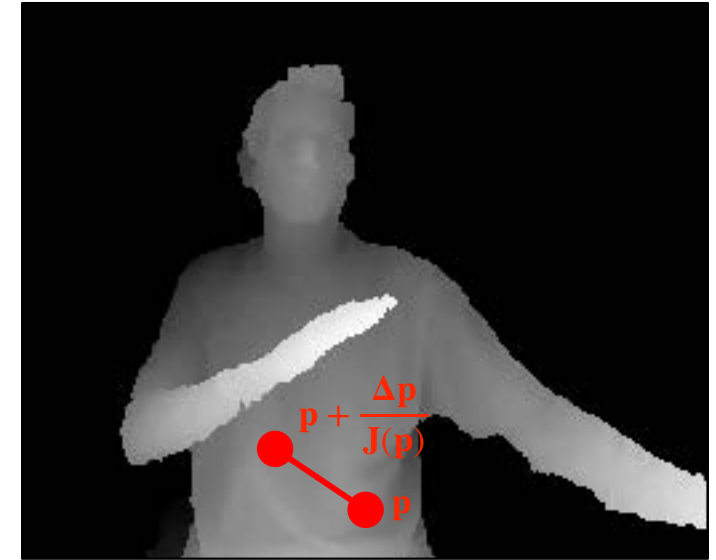
# Body Part Recognition



Input depth image



Training labelled data



Visual features

Visual feature:

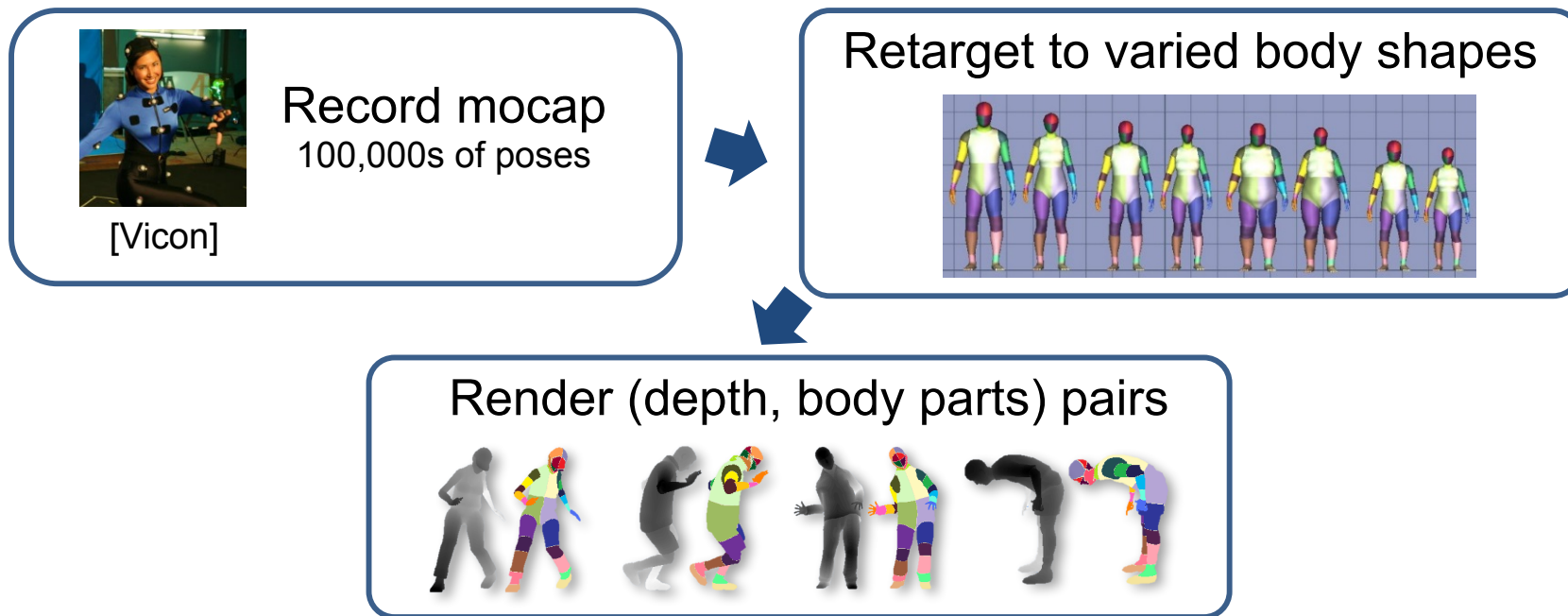
$$x(\mathbf{p}, \Delta\mathbf{p}) = J(\mathbf{p}) - J\left(\mathbf{p} + \frac{\Delta\mathbf{p}}{J(\mathbf{p})}\right)$$

Weak classifier:

$$h(\mathbf{p}, \Delta\mathbf{p}, \tau) = x(\mathbf{p}, \Delta\mathbf{p}) - \tau$$

- Very fast to compute.
- Real-time performance

# Synthetic Training Data



Train invariance to:





# Influence of Tree Depth

Input depth



Ground truth parts



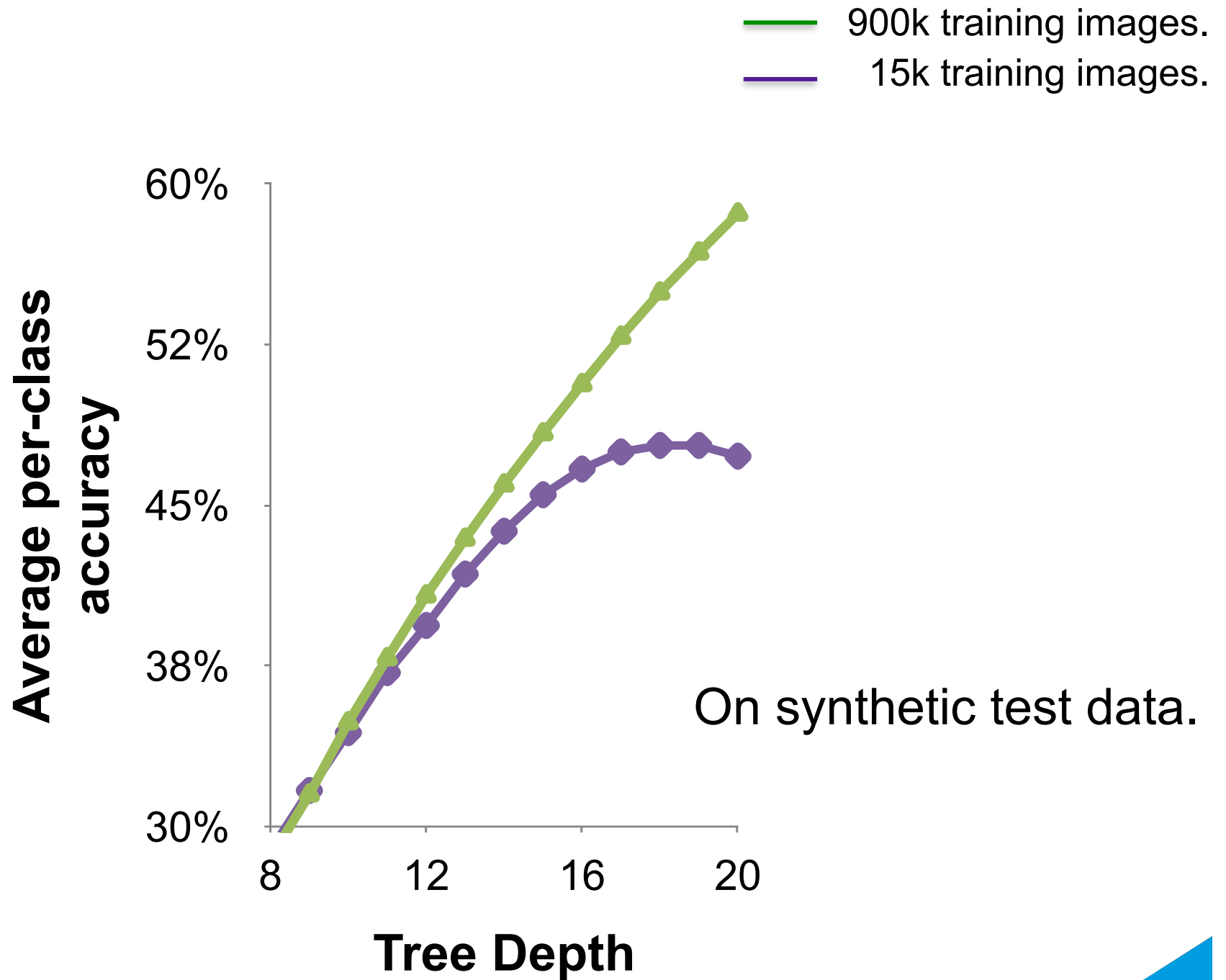
Inferred parts



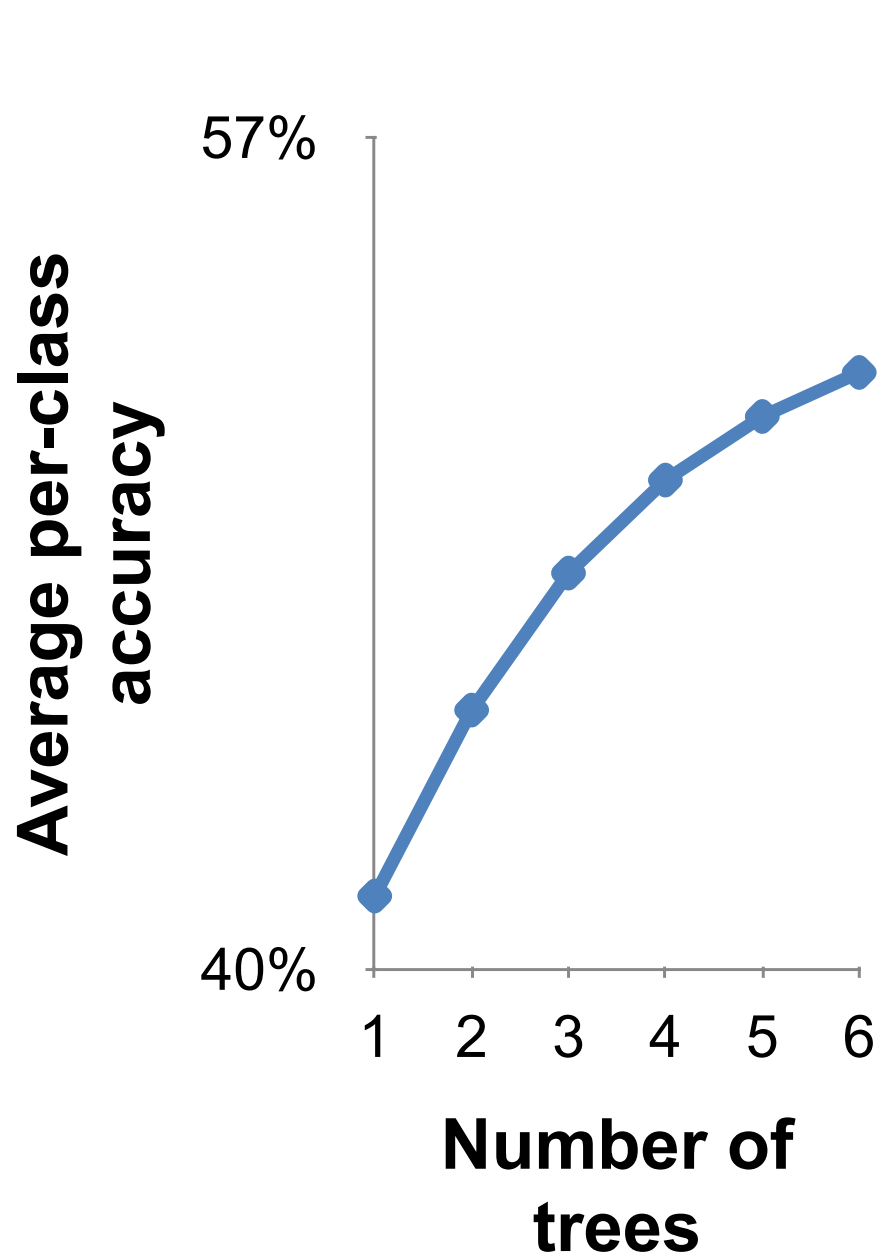
depth 18



# Choosing the Tree Depth



# Choosing the Number of Trees



ground truth



inferred body parts (most likely)

1 tree



3 trees



6 trees



# Result



Input depth image with background removed.



Inferred body parts posterior  
 $p(c|\mathbf{v})$

# Decision Forests in Short

- They make it comparatively easy to interpret what is happening.
- Their behavior is easy to modify.
- They can be trained using moderate amounts of data.

—> Very useful in many practical applications.