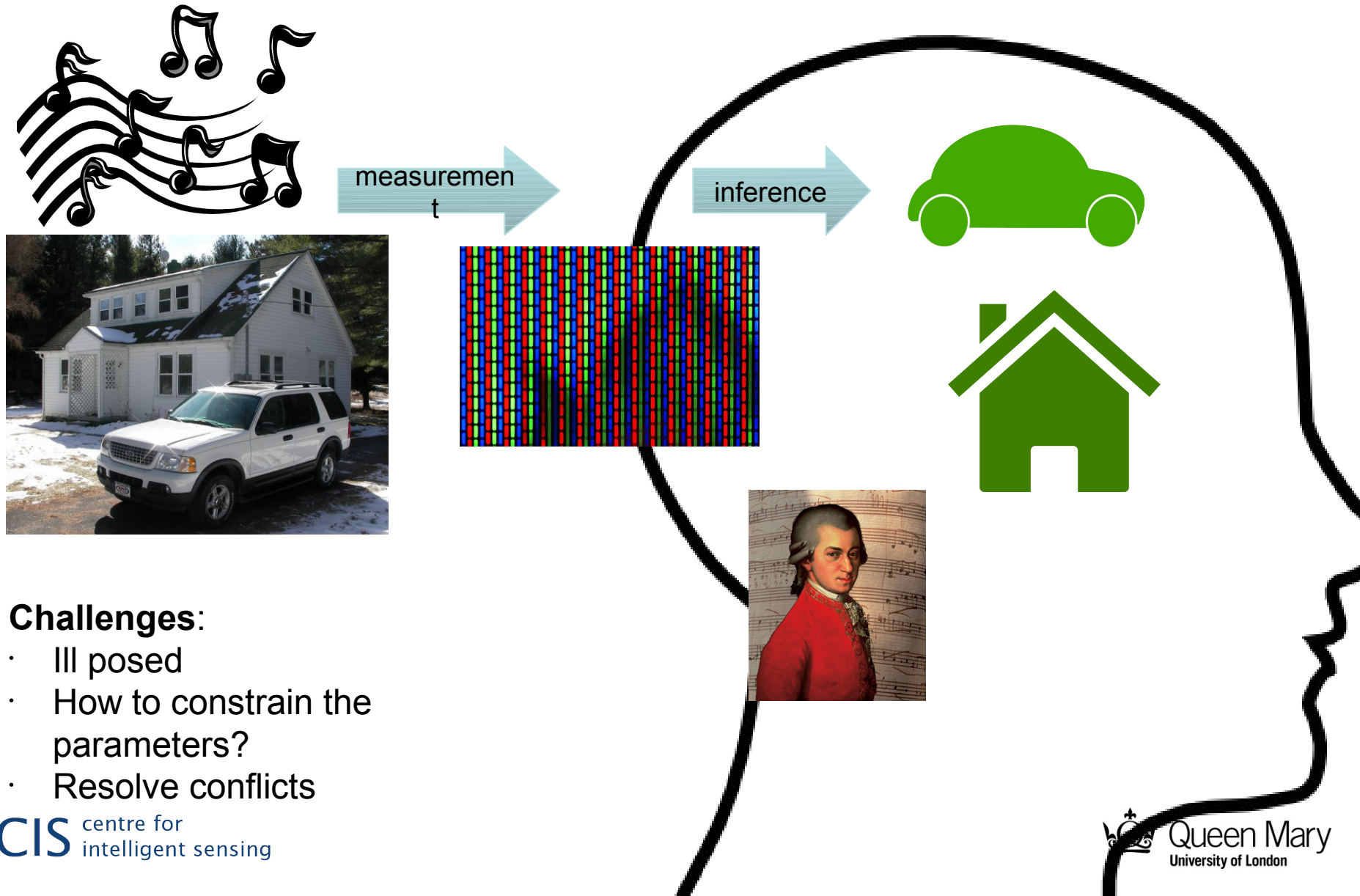

Probabilistic Machine Learning Models for Intelligent Sensing: Computer Vision and Beyond

Dr. Timothy Hospedales

Centre for Intelligent
Sensing

Queen Mary University of
London

Intelligent Sensing?



Challenges:

- Ill posed
- How to constrain the parameters?
- Resolve conflicts

Intelligent Sensing?

Intelligent Sensing

- How to resolve different sources of data?
- How to solve ill-posed problems?
- How to not start from scratch every time?
- How to make the right measurements?

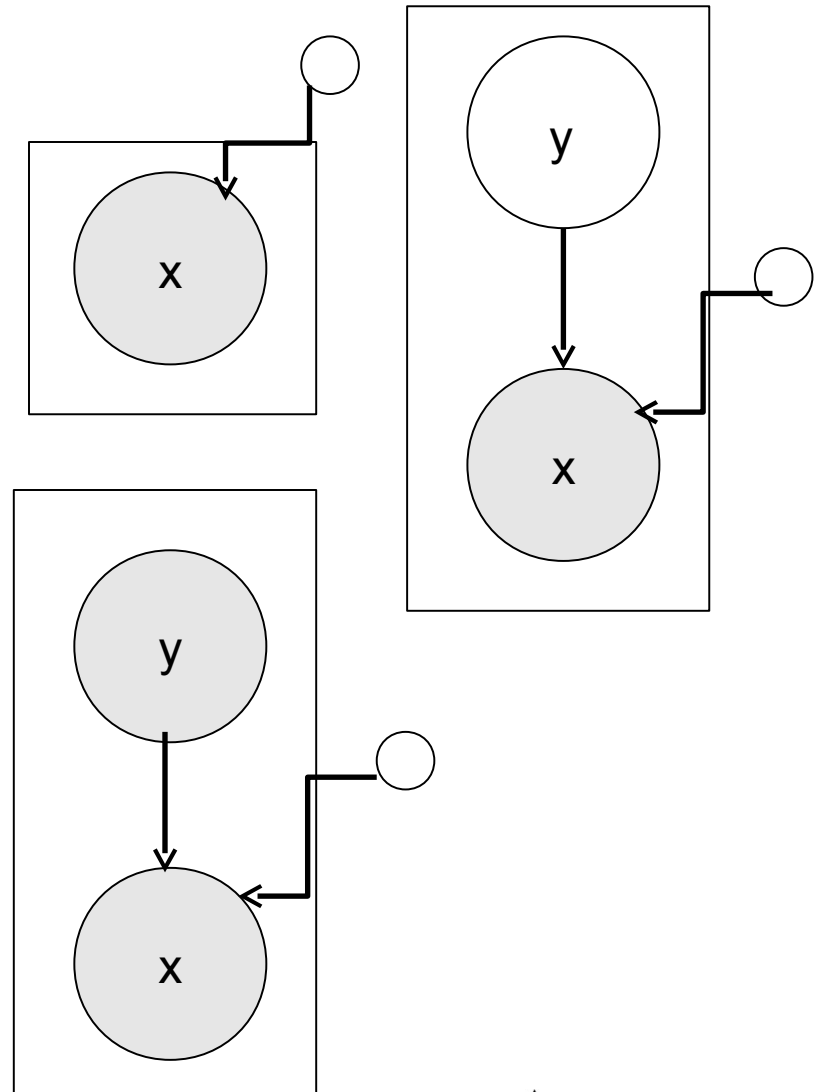
Outline

Intelligent Sensing

- **Data Fusion**
- **Data** and **annotation** efficiency
 - Unsupervised
 - Weakly-supervised
 - Semi-supervised
 - Multi-label/multi-instance
 - Zero-shot learning
- **Observation** efficiency
 - Active learning

Notation

- Unsupervised Learning
 - Observe $\{x\}$, model $p(x)$
 - Observe $\{x\}$, model $p(x,y)$
- Supervised Learning
 - Observe $\{x,y\}$, model $p(x,y)$



Machine Learning

Classic Problems:

- Inference
- Marginal Likelihood
- ML Learning,
 - Density Estimation
- EM Learning
- Model Selection

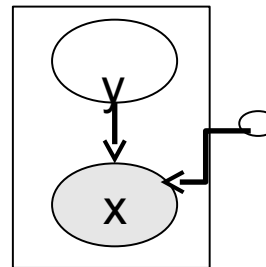
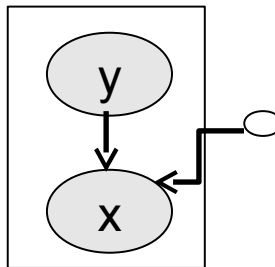
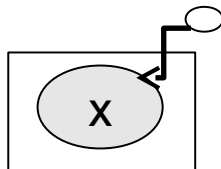
$$p(y | x) = p(x, y) / p(x)$$

$$p(x) = \int p(x, y) dy$$

$$\hat{\theta} = \operatorname{argmax}_{\theta} p(X | \theta)$$

$$\hat{\theta} = \operatorname{argmax}_{\theta} \int p(X, Y | \theta) dY$$

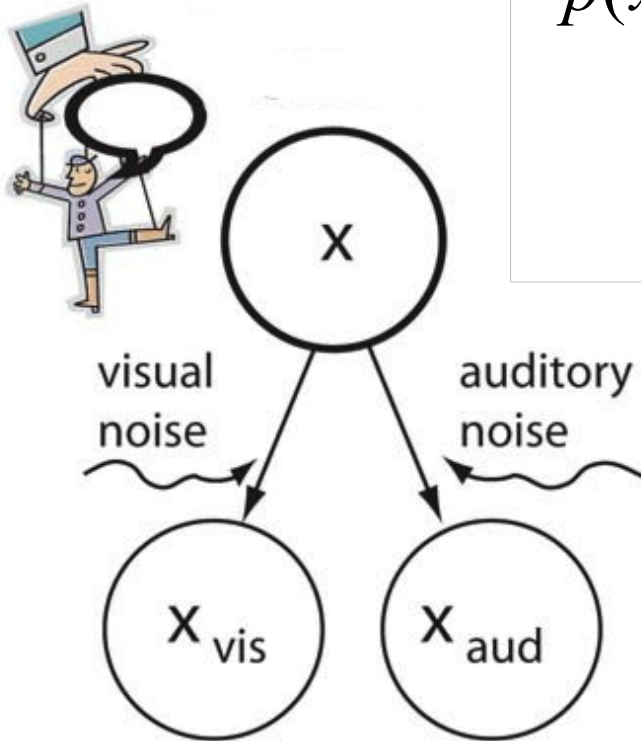
$$M = \operatorname{argmax}_{M} \int p(X, Y, \theta | M) p(M) dY d\theta$$



-
- Fusion and Data Association

Multisensory perception

Fusing Multiple Data Sources



$$\begin{aligned} p(x | x_a, x_v) &\propto p(x_a, x_v | x)p(x) \\ &= p(x_a | x)p(x_v | x)p(x) \end{aligned}$$

Optimal fusion depends on reliability of each modality

But how do you know the reliability of your senses?

Fusing Multiple Data Sources

Aim: Given an audio-visual stream

- Learn the user's appearance & sound

- Learn the microphone and camera characteristics

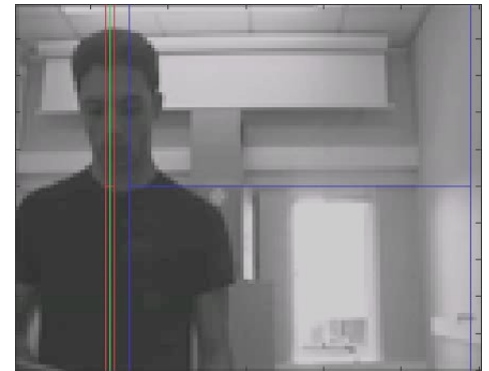
- Fuse appearance and sound info for optimal localization

Challenges:

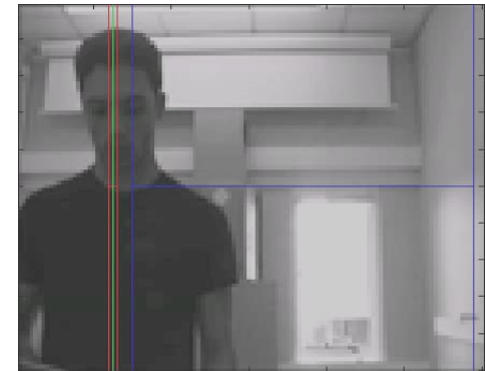
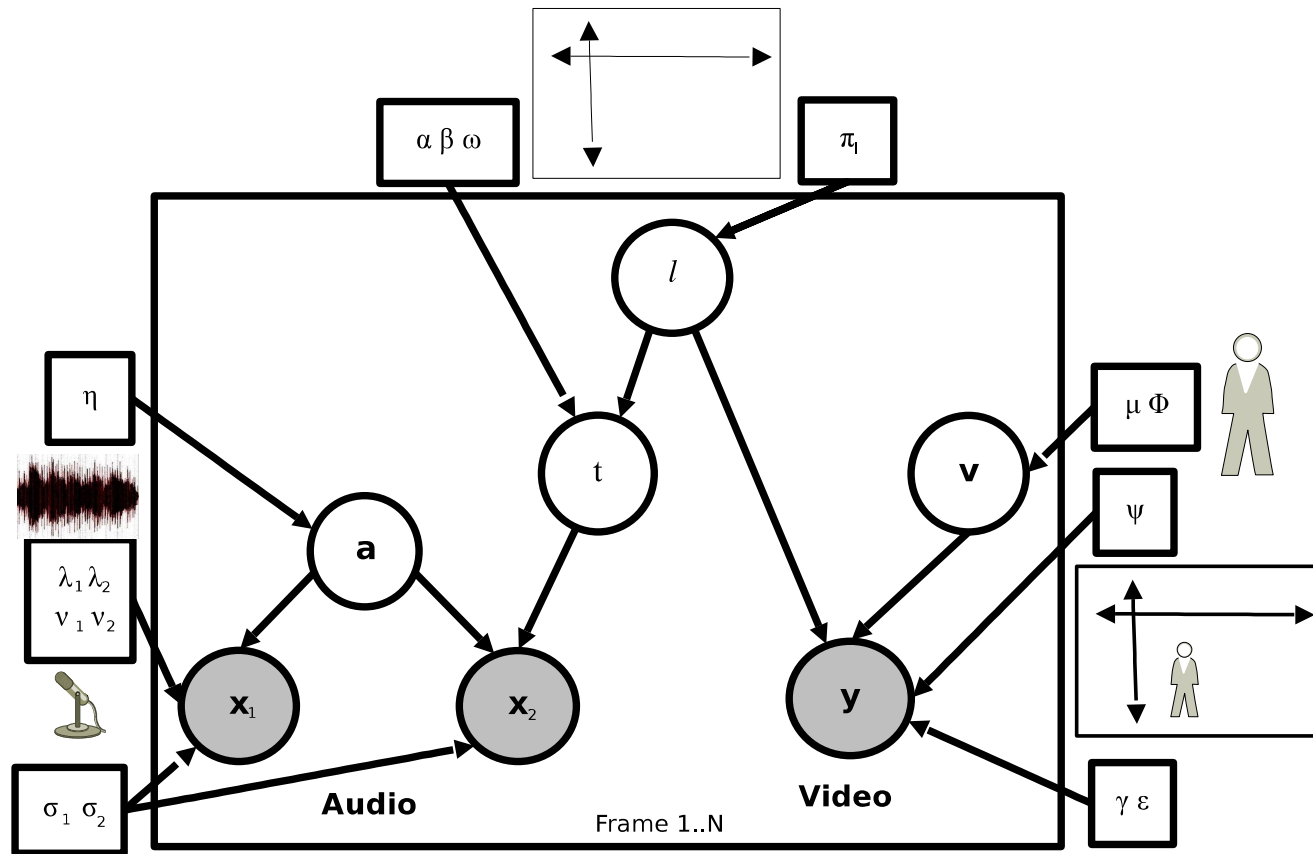
- No supervision

- (No background subtraction)

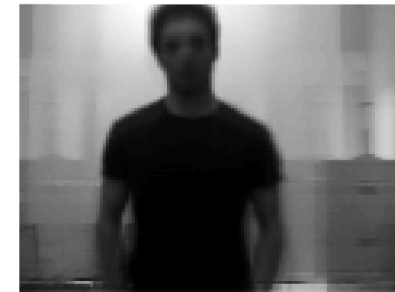
- Real-time inference



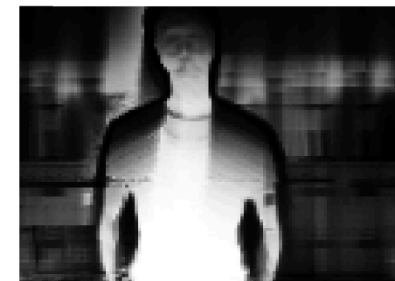
Fusing Multiple Data Sources



(f) Learned Template Mean

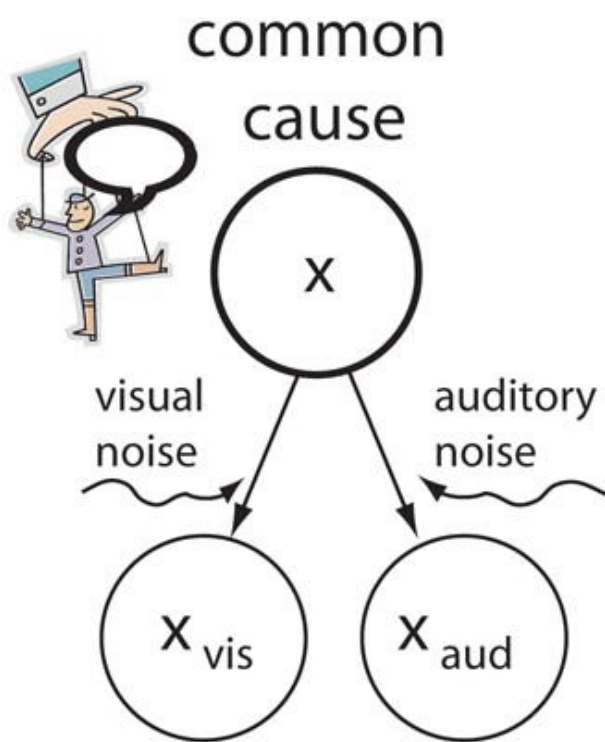


(g) Learned Template Precision

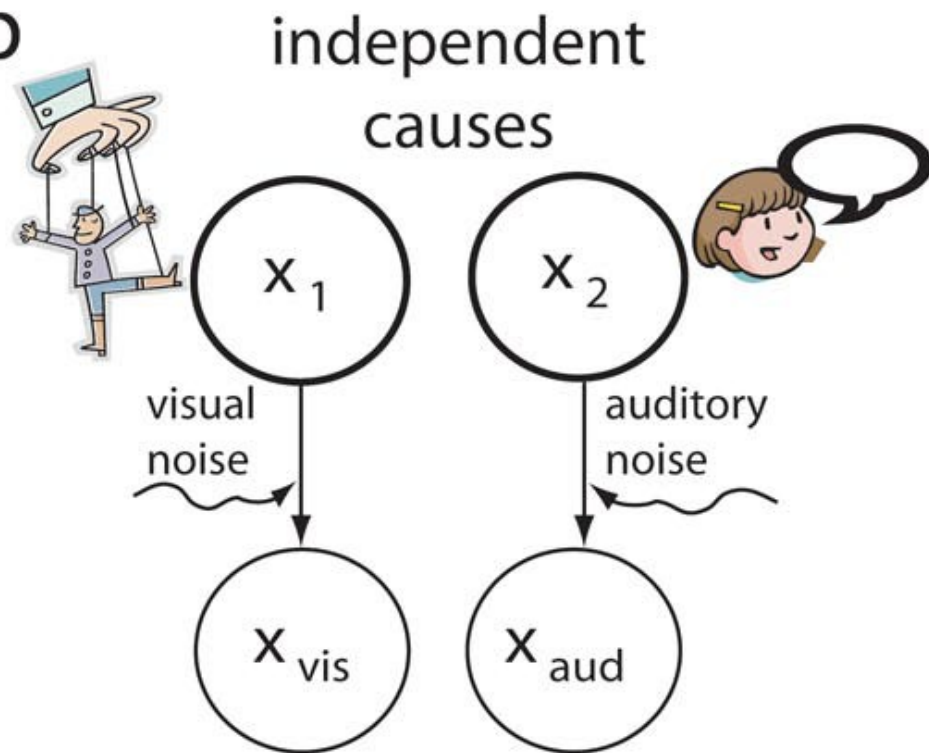


Fusion without correspondence?

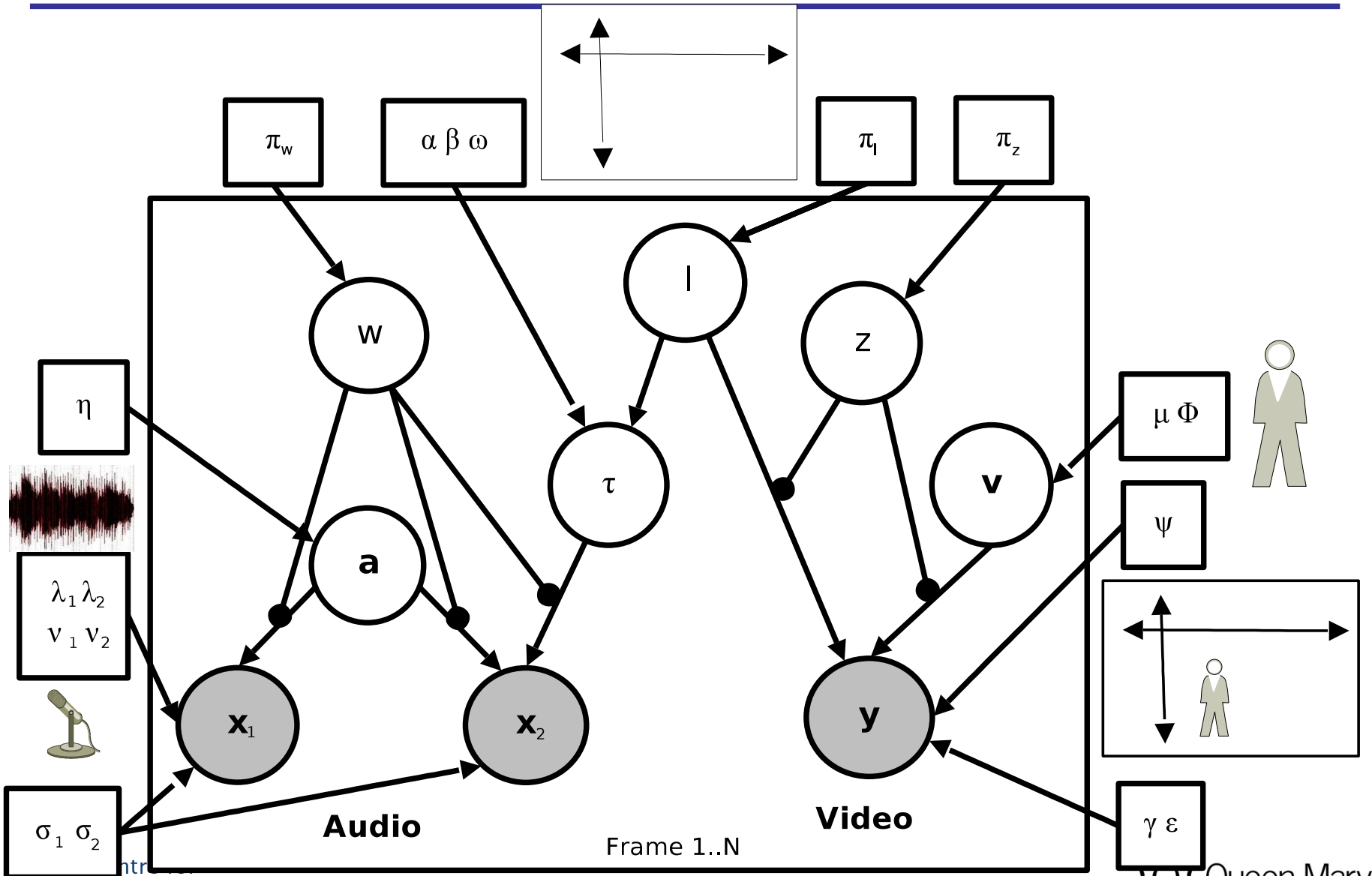
a



b



Multiple Data Sources?



Fusion without correspondence?

Unsupervised AV Scene Understanding: Who Said What, Where and When?



Multitarget Tracking

User 1 not visible.



User 1 silent.

User 2 not visible.



User 2 silent.

EM Learning
Model Selection

Appearance Learning

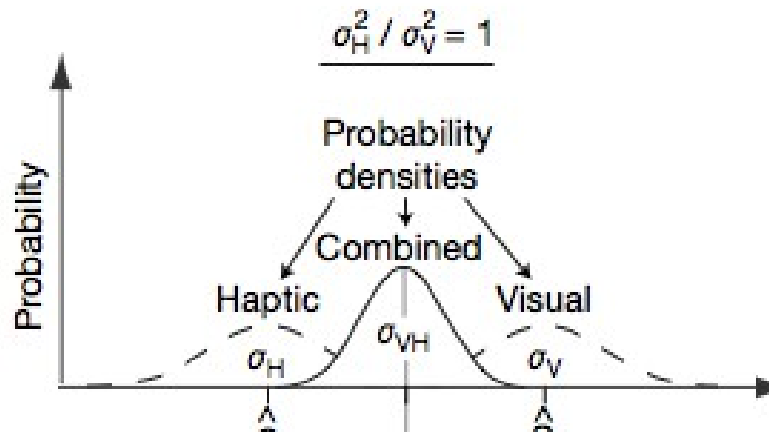
AV Data Association

-
- Abnormality Detection

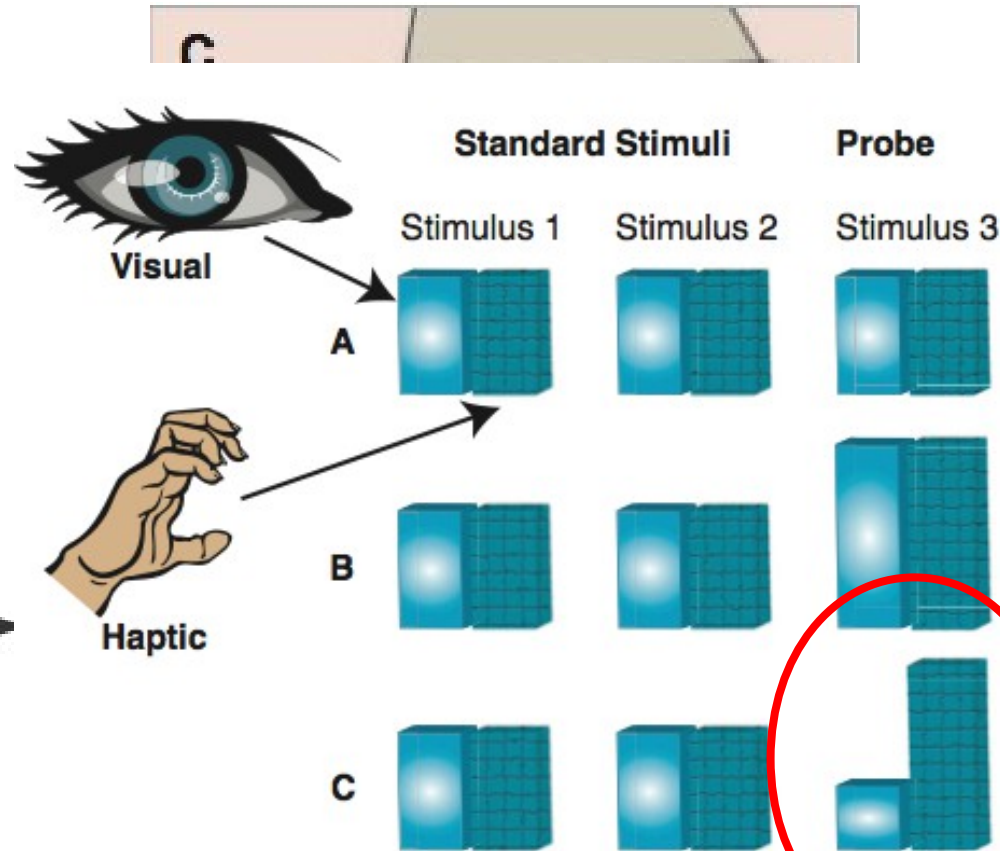
Multisensory perception

Human Multisensory Oddity detection

Optimal Fusion

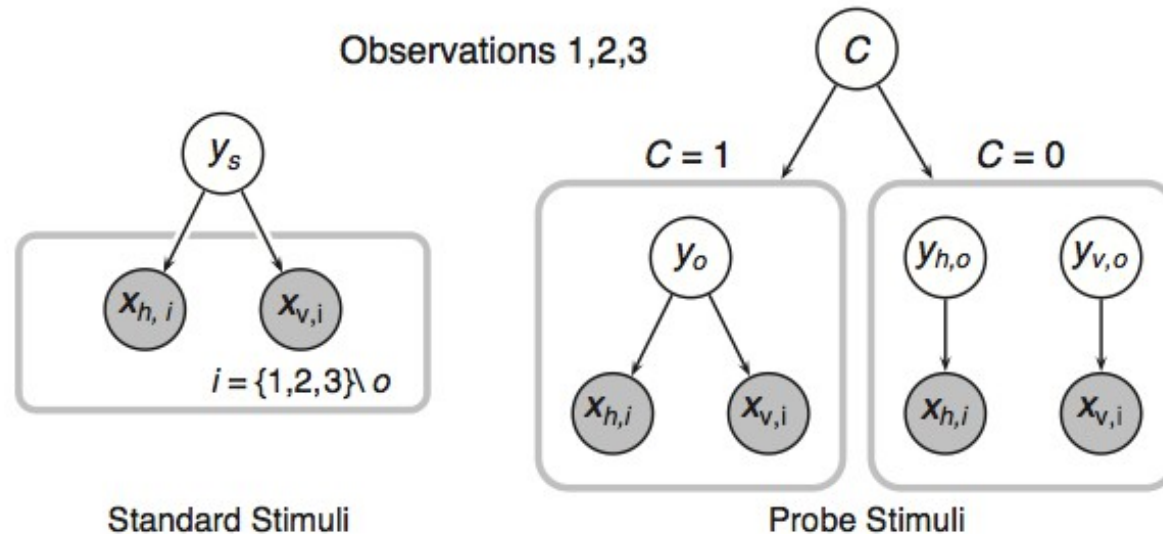


$$\hat{S}_c = w_1 \hat{S}_1 + w_2 \hat{S}_2$$

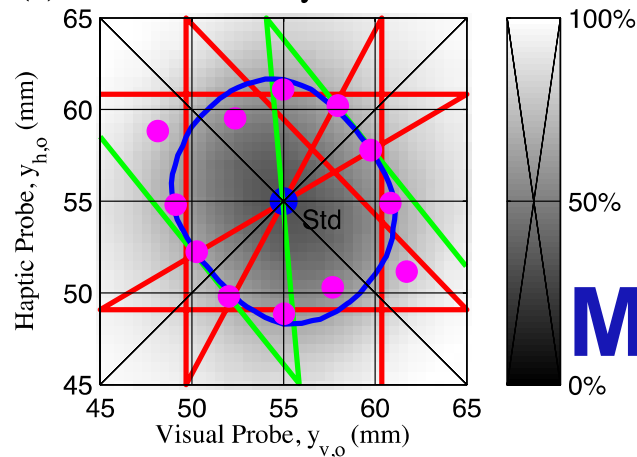


People (eventually)
notice this

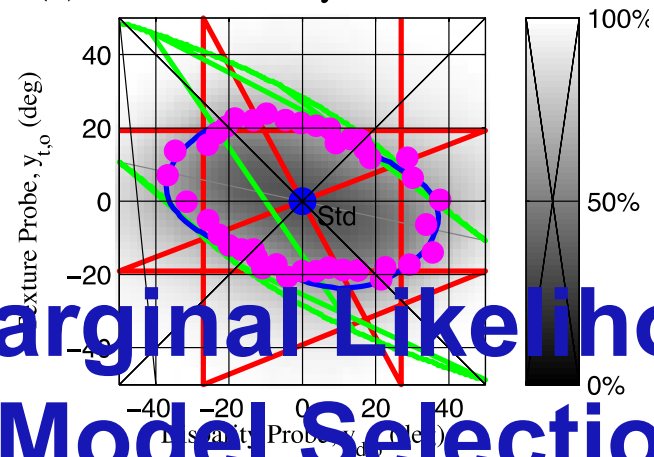
Human Multisensory Oddity detection



(a) Across Modality Detection Rate



(b) Within Modality Detection Rate



**Marginal Likelihood
Model Selection**



Data Fusion

- Optimally integrating multiple sensors
- Resolving multi-sensor data association

Outline

Intelligent Sensing

- Data Fusion
- **Data and annotation efficiency**
 - Unsupervised
 - Weakly-supervised
 - Semi-supervised
 - Multi-label/multi-instance
 - Zero-shot learning
- Observation efficiency
 - Active learning

-
- Video Surveillance: Anomaly Detection & Clustering
- ## Unsupervised learning

Unsupervised Learning / Surveillance

- Aim: Given a video stream
 - Get domain knowledge by learning about activities
 - Detect abnormal behaviors as outliers against the model
- Challenges:
 - No tracking
 - No camera calibration
 - No supervision
 - Complex behaviors
 - Real-time

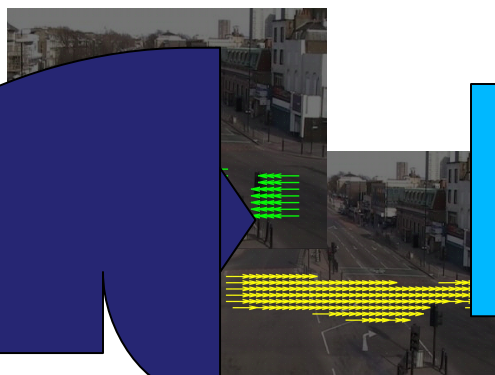


Unsupervised / Surveillance: MCTM

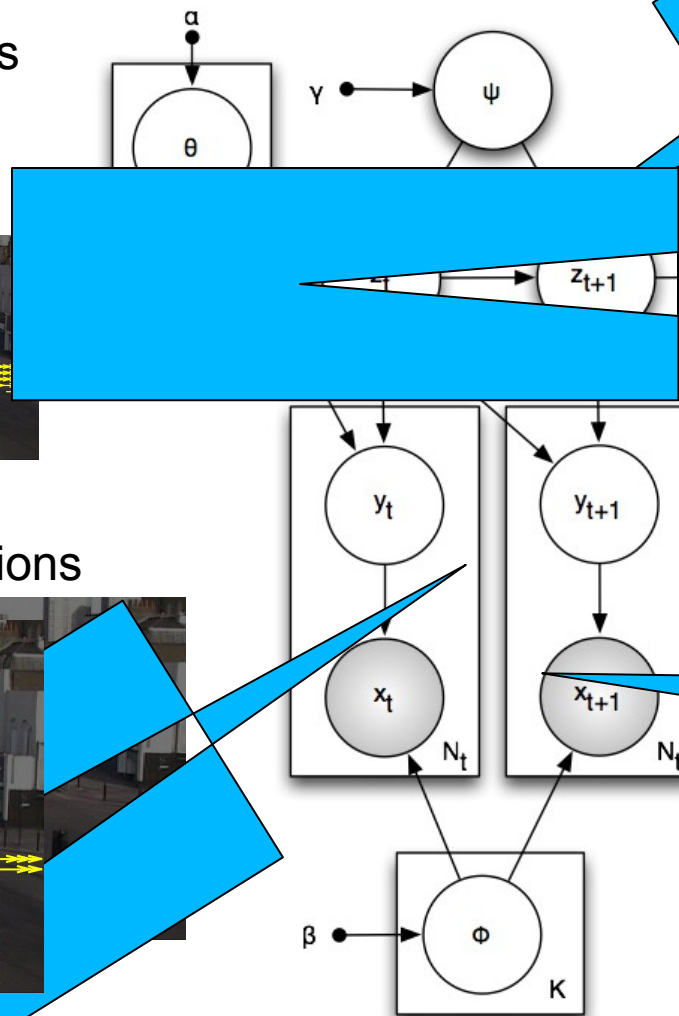
New generative model...

“Markov Clustering Topic Model”

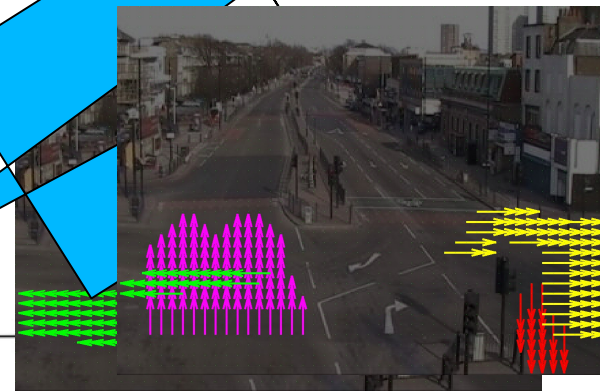
Learned Dynamics



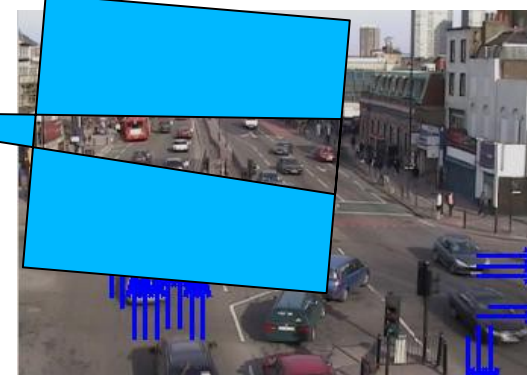
Learned Actions



Learned Behaviors



Input Flow Vectors



Unsupervised / Surveillance: MCTM Learning



Unsupervised / Surveillance: MCTM: Abnormality



$$p(x) = \int p(x, y) dy$$

$$p(y|x)$$

Unsupervised Learning



- What if you have seen a few known (**but rare**) behaviors you want to detect?
 - E.g., surveillance
- These may also be **subtle**

Weakly Supervised Learning / Surveillance

- Aim: Given a video stream
 - Learn a detector for a rare and subtle behavior
- Challenge
 - Use only **weak annotation**
 - **Sparse training data**
 - No tracking
 - Real-time
- View as hard multi-instance learning (MIL) problem

WSL / Surveillance: Rare Events

Example Challenge

1 Example Each



100 Examples



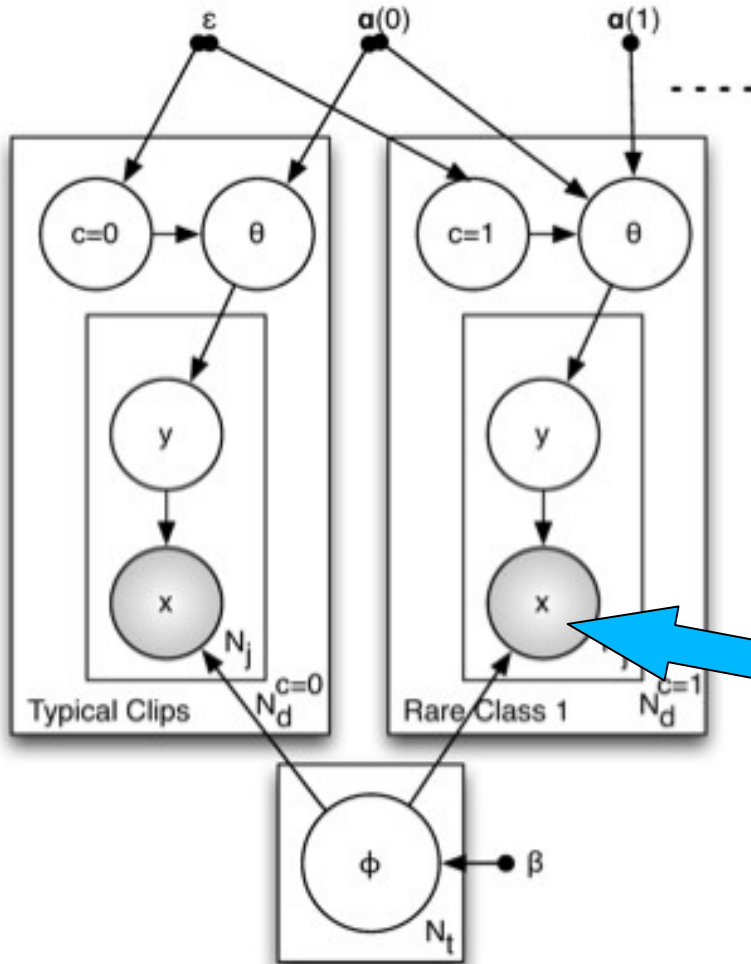
WSL / Surveillance: WSJTM

Classify: Compute $p(C|X)$

- Bayesian Model Selection
- Variational Importance Sampler

Locate: Infer $p(Y|X, C)$

- Gibbs



WSL / Surveillance: WSJTM: Results

WSJTM Classifier: Trained with Weak and Sparse Labels



EM Learning
Model Selection

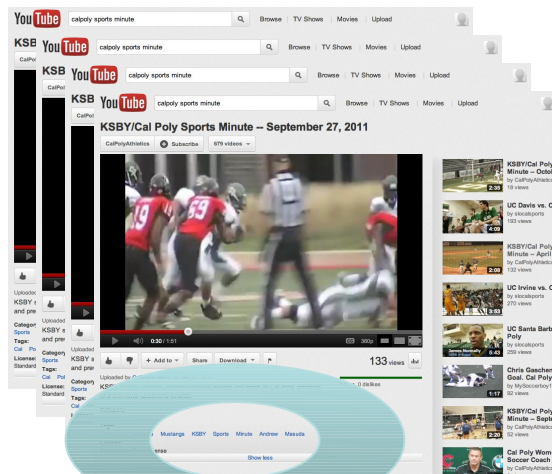
Weakly Supervised Learning



- What if you need more than one label per instance?
 - E.g., multi-media indexing.
- The weak supervision problem gets harder...

Weakly Supervised Multi Label: Tagging

- Aim: Learn video annotation model from tags
 - Online video databases: Sports news



Tags

Learning

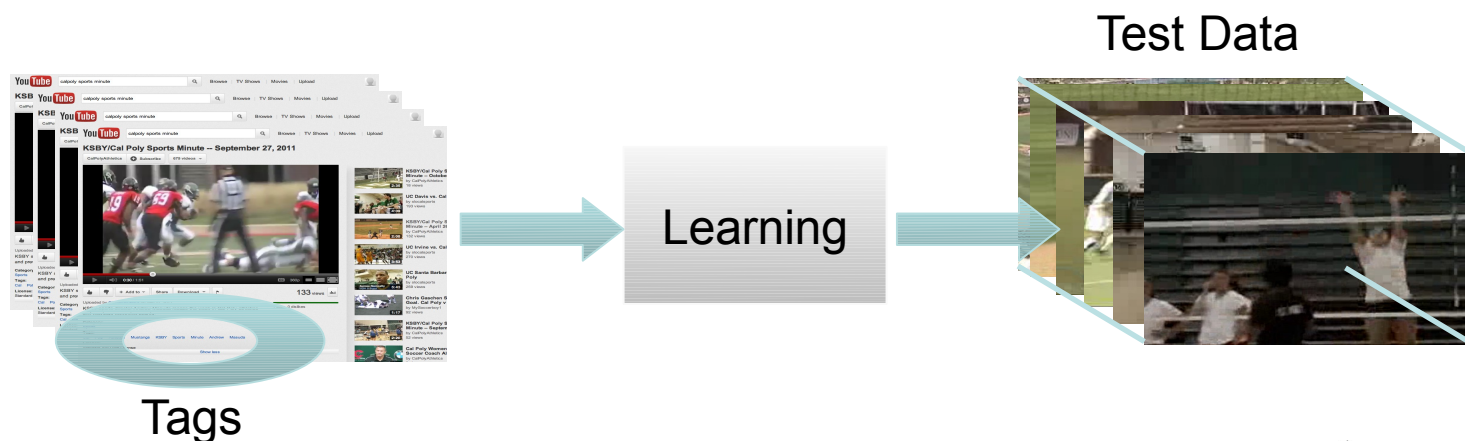
Test Data



Annotation:
Baseball, Soccer, Volleyball

Weakly Supervised Multi Label: Tagging

- Challenges:
 - Weak annotation
 - Multiple labels per instance
 - Huge intra-class variability



Weakly Supervised Multi Label: Results



Weakly Supervised Multi Label: Insight

- Foreground Topics
 - Shared, e.g., running
 - Specific, e.g., pitching



EM Learning Inference

Weakly Supervised / Multi Label Learning



- Can we learn anything **with no new data at all?**
 - (Classification)

Zero-Shot Attribute Learning / Tag & Classify

- Aim: Given a set of Tags & Classes
 - Learn how to tag
 - Learn how to relate tags to classes
 - **Zero-shot learning** from tag description

Zero-Shot Learning (Look Mum! No Data!)



Stripes, Herbivore, Tail, Claws
→ Zebra



Lion!

Lion:= Stripes, Herbivore, Tail, Claws

Zero-Shot Learning (Look Mum! No Data!)



Stripes, Herbivore, Tail, Claws
→ Zebra



Lion!

Lion := Stripes, Herbivore, Tail, Claws



Candles, Cake, Clapping, Dancing
→ Birthday Party

Wedding Dance := Candles, Cake,
Clapping, Dancing

Zero-Shot Learning (Look Mum! No Data!)



Stripes, Herbivore, Tail, Claws
→ Zebra

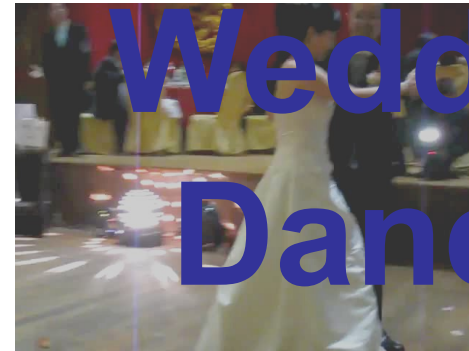


Lion!

Lion:= Stripes, Herbivore, Tail, Claws



Candles, Cake, Clapping, Dancing
→ Birthday Party



Wedding Dance!

Wedding Dance:= Candles, Cake, Clapping, Dancing

Latent Attribute Learning

- Aim: Given a set of Tags & Classes
 - Learn how to tag
 - Learn how to relate tags to classes
 - Zero-shot learning from tag description
- Challenge: Reduce human effort
 - Avoid annotating every attribute on every training image
 - Avoid specifying every attribute on every new class

Latent Attribute Learning

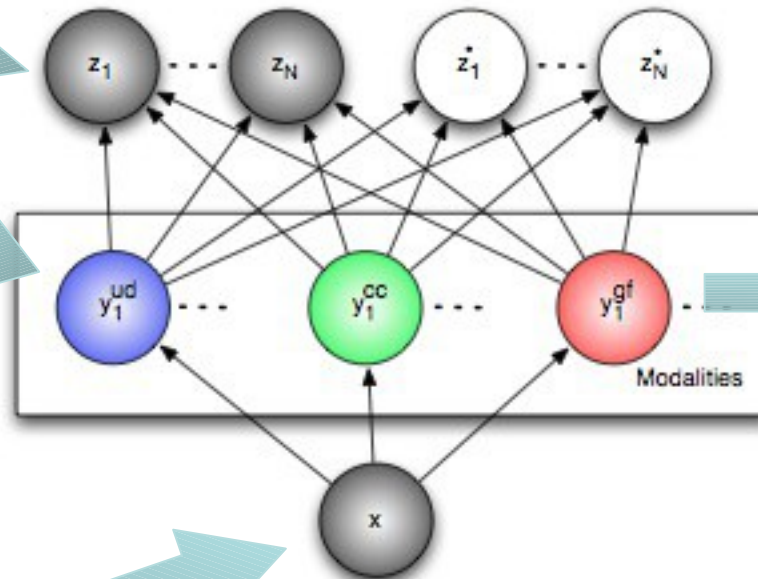
Classes

Lion
Zebra

Attributes

Stripes, Herbivore,
Tail, Claws
Stripes, Herbivore,
Tail, Claws

Image/Video



Latent
Attributes

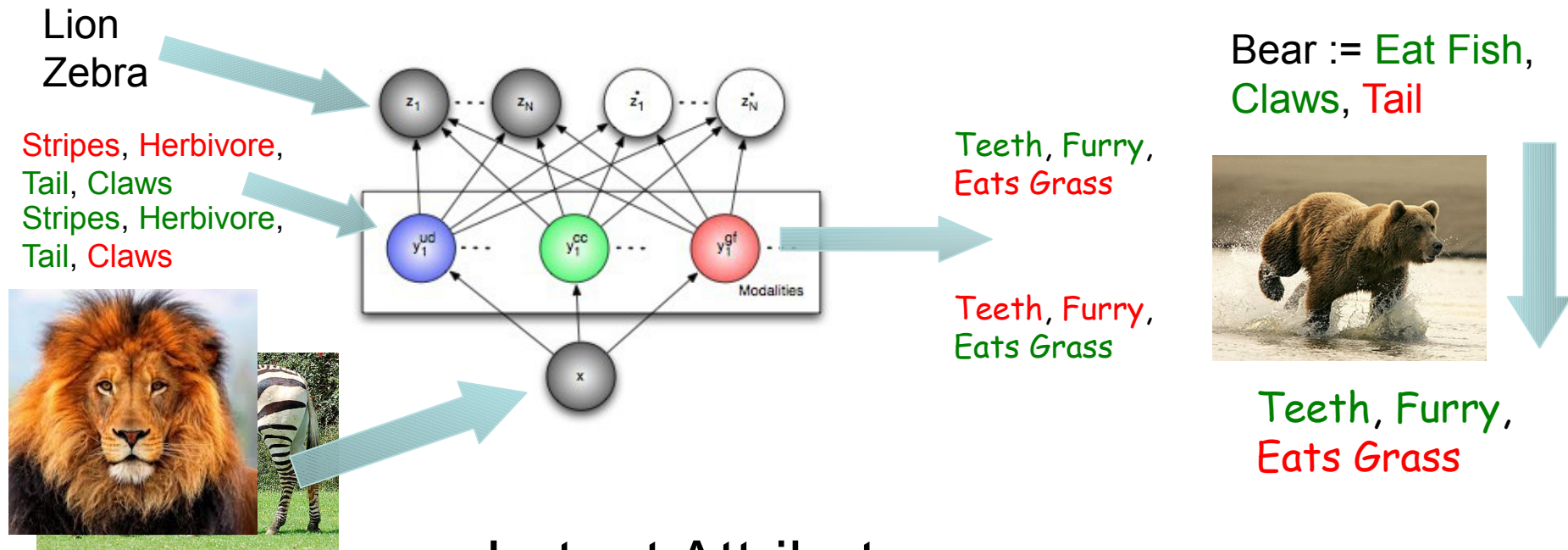


Teeth, Furry,
Eats Grass

Teeth, Furry,
Eats Grass



Latent Attribute Learning



Latent Attributes:

- Less annotation work
- Increased Classification Accuracy
 - Conventional
 - Zero-shot

Zero-shot learning



- We can also use attributes for **sparse data** learning
 - E.g., in **re-identification**.

More on attributes: Re-identification





More on attributes: Re-identification

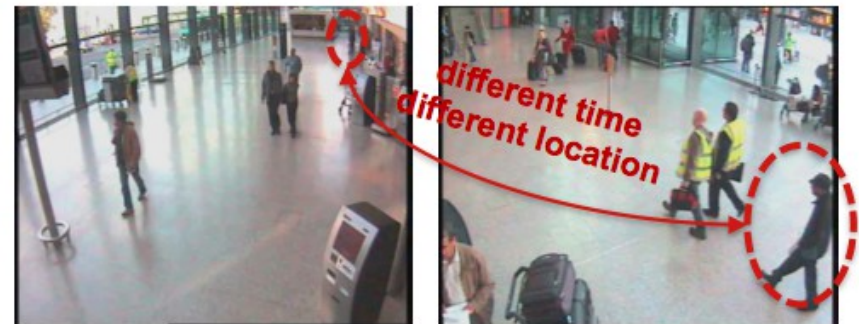
Aim: Re-identify across time, space and view

Solution?

- ~~Learn a recognizer~~

Challenge:

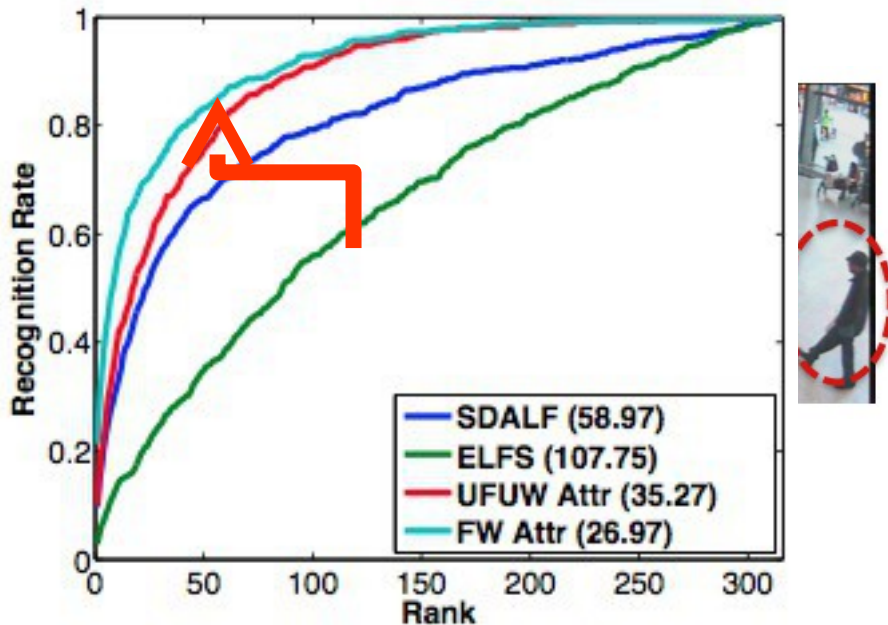
- Statistical Insignificance
- (One-shot learning)
- Huge intra-class variability



More on attributes: Re-identification

Aim: Re-identify across time, space and view

Solution: Attribute Transfer



- Target Person

- ✓ Hat
- ✗ Jeans
- ✓ Male
- ✓ Coat
- ✗ Skirt
- ✗ Tie
- ✗ Sandals
- ✗ Shorts

- Leverage lifetime of (attribute) experience

Sparse-data and re-identification



What if we have a lot of data but not much supervision?

- Active Learning

Outline

Intelligent Sensing

- Data Fusion
- Data and annotation efficiency
 - Unsupervised
 - Weakly-supervised
 - Semi-supervised
 - Multi-label/multi-instance
 - Zero-shot learning
- **Observation efficiency**
 - Active learning

Active Learning

Aim: Make query selection optimal

→ Minimize human annotation effort for given outcome



Challenge:
Which Example Will Help Me The Most?

Active Learning & Discovery

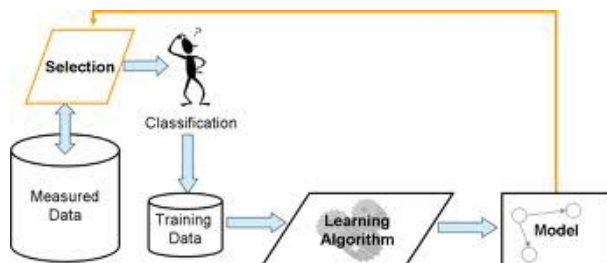
Aim: Make query selection optimal

→ Minimize human annotation effort for given outcome

Challenge:

How to do in a new domain with unknown class space?

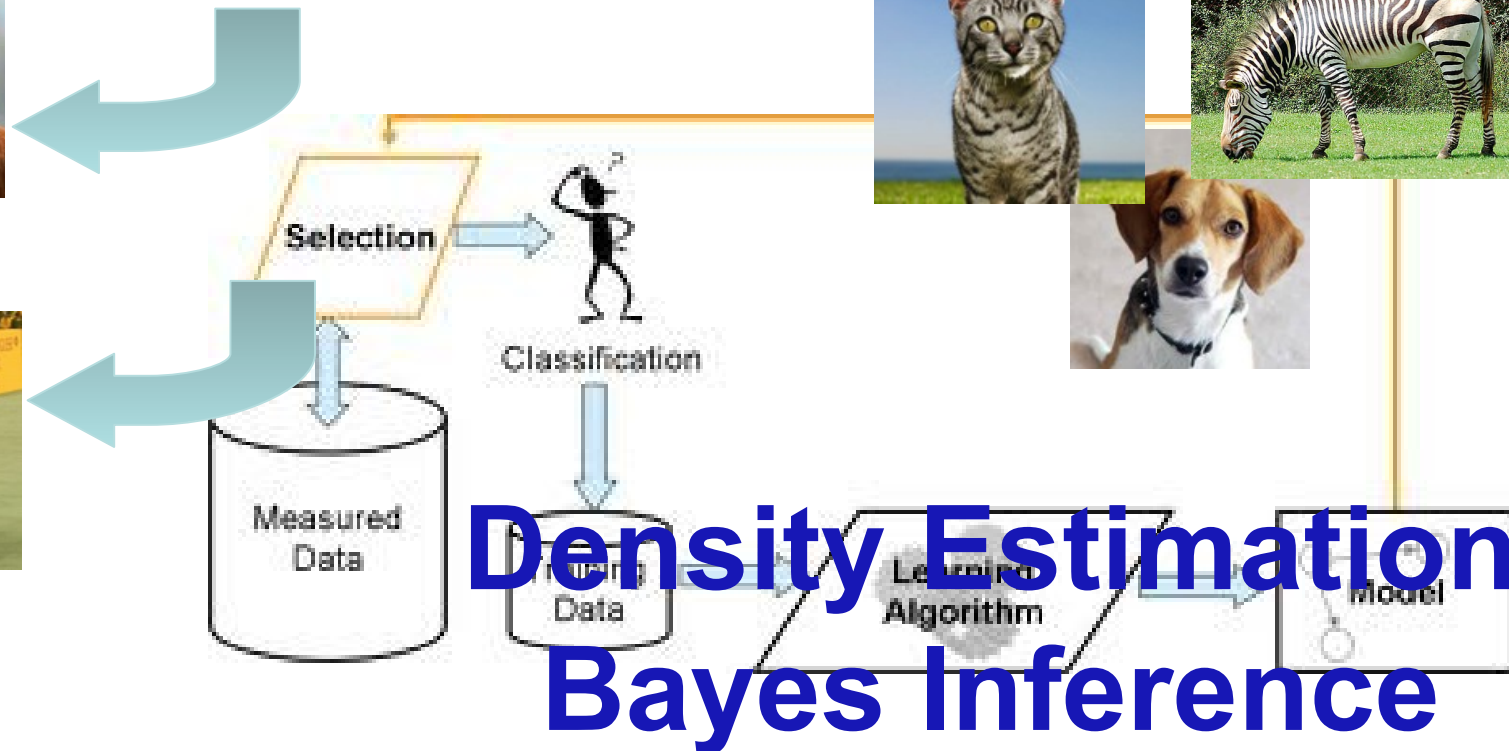
Tractability?



Active Learning & Discovery

Solution:

- Bayesian non-parametrics
- Incremental Computation



Active Learning & Discovery



Take Homes

- Be aware of the underlying ML of your intelligent sensing problem
 - Then you can find good techniques
- Separation of:
 - features, objectives, model/representation, optimizers
- Think of your data and annotation constraints
 - Can they be reduced to make your model more useful?
 - How does it depend on the strength of your annotation?
 - Could your model do better with more data (but same annotation?)
 - Would finding the right annotations help?