

(1) Compositional Models: Complexity of Representation and Inference

Note Title

11/24/2013

Beyond Feedforward Models:

"something about generative models"

(I) General Comments: Data Driven Models and Representations.
Feed-forward and Feedback Architectures.

(II) Unsupervised Learning of Representations and Hierarchies.

(III) Complexity Results on Representations and Inference.

(2) (I) General Comments: Data Driven Models

Note Title

11/23/2013

In the last 10 years there has been considerable success using data driven models.

- Face Detection. (Viola & Jones)
- Text Detection (Cheu & Yuille)
- Pascal Challenge. DPM's. (Felzenszwalb, Ramanan, McAllester)
Object Detection. L. Zhu et al.
- Image Net. DBN's Krizhevsky et al.

Question: Are datasets big and representative enough?
Will results scale? Yes and No.

(3) (I) General Comments : Representations

Vision addresses an enormous range of tasks:

single object: detect in clutter and with occlusion, detect parts and boundaries, estimate 3D structure, reason about pose.

multiple objects: positions relative to each other, occlusion relations, social interactions, etc.

Scene structure: ground plane, Manhattan World structure (if opprop), positions of objects in scene, background stuff (sky, water), surfaces, geometry, motion.

Visual system: needs to compute rich representations from images.

(4) (I) General Comments: Feedforward and Feedback

• Classic Feedforward Theories:

Marr: Primal Sketch \rightarrow 2.5D Sketch \rightarrow 3D Rep

• Hierarchical Feedforward Models: (Invariances)
Fukushima, HMax, Deep Belief Networks, M-Theory

Feedforward and Feedback: Generative Models

• Analysis by Synthesis - Mumford & Grenander.

• DDM (MC) - Tu & Zhu, Tu, Chen, Yuille, Zhu
feedforward proposals validated by feedback models.

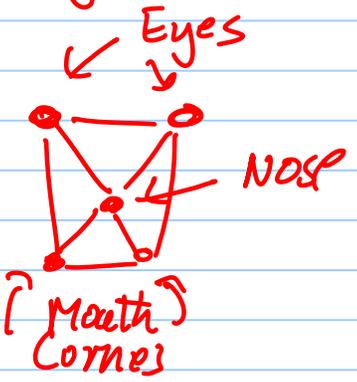
But do generative models require feedback?

Ullman:
Streams
Counter-Streams

(5) (11) Unsupervised Learning of Representations and Hierarchies

Starting Point : Pictorial Structures. (Fischler & Everslager) 1973

Model Object (Face) in terms of parts



$G = (V, E)$ Graph

state variable $\{x_\mu : \mu \in V\}$
 μ nodes, E edges

energy

$$E[\langle x_\mu \rangle] = \sum_{\mu \in V} \lambda_\mu \phi(x_\mu) + \sum_{(\mu, \nu) \in E} \lambda_{\mu\nu} \psi(x_\mu, x_\nu)$$

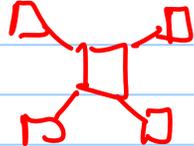
data term. spatial relations

Deformable Template:

(G) (II) Unsupervised Learning of Representation and Hierarchies.

These models are very successful — often by putting many features into the data term.

Examples.

- von der Malsburg — Face Detection "Neurcl".
Wiscott., Neven
- Coughlan, Sutte — Hands (Dynamic Programming)
- Felzenszwalb, Ramanan, McAllester — Star Model, DPM
(Pascal Detection) 

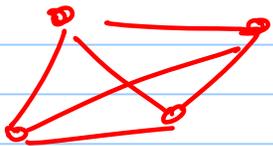
(7) (II) Unsupervised Learning of Representations and Hierarchies

Unsupervised Learning

Constellation Models.

Caltech.

- Weber, Perona
- Fergus et al.

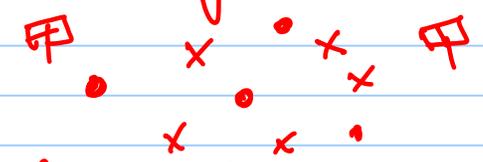


Fully connected?

Flat Models → no composition

(2) (II) Unsupervised Learning of Representation and Hierarchies

L. Zhu et al.: Learn Mixtures of Models
Data with clutter

• Take Image \rightarrow Extract and Represent Interest Points (IPs)
 Interest Point types. $\#I \cdot X$
(~ 200).

• Task - Learn a Generative Model.

Don't know:

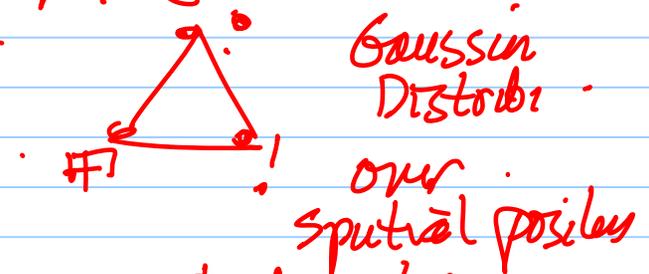
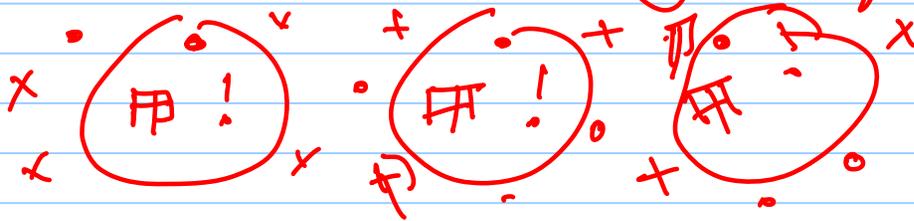
- (i) how many objects
- (ii) how many IP's in an object
- (iii) which IP's are object or background.

(a) (II) Unsupervised Learning of Representations and Hierarchies

Strategy: Greedy search over space of models

Initialize: Default model all IP's generated independently (i.i.d)

Cluster: Identify frequent trip lets



"Suspicious coincides" Barlow
"meaningful alignments" Morel.
'.....' Ullman.

Note: Gaussian over internal angles
invariant to rotation / scale.

(17) (II) Unsupervised Learning of Representations and Hierarchies

Grow Model by Adding Triplets:

Current Model:

IP's generated by i.i.d. background or triplet.

Image. N background point B
+ $3M$ points generated by \triangle Triplet Model.

B \triangle
 \downarrow \downarrow Rep by $B + \triangle$

Grow to $B + \triangle$ or $B + \triangle + \nabla$
mono complex object ∇ new object

(11) Unsupervised Learning of Representation and Hierarchies

Better Encoding of Data:

• Cost of Encoding Data by B only Default!

• Cost of Encoding Data by $B + \Delta$

$$\sum \log P(\dots) + T_{\text{inferred}}$$

Model Selection?

Stop when adding "new object"
or growing b_i

(12) (II) Unsupervised Learning of Representations and Objects

Learn Representations → Unsupervised.
→ I.P.'s only (less interesting)

Learning in Cocktail Party.

Unknown no. of speakers + Background Noise.

Can do non-trivial vision tasks

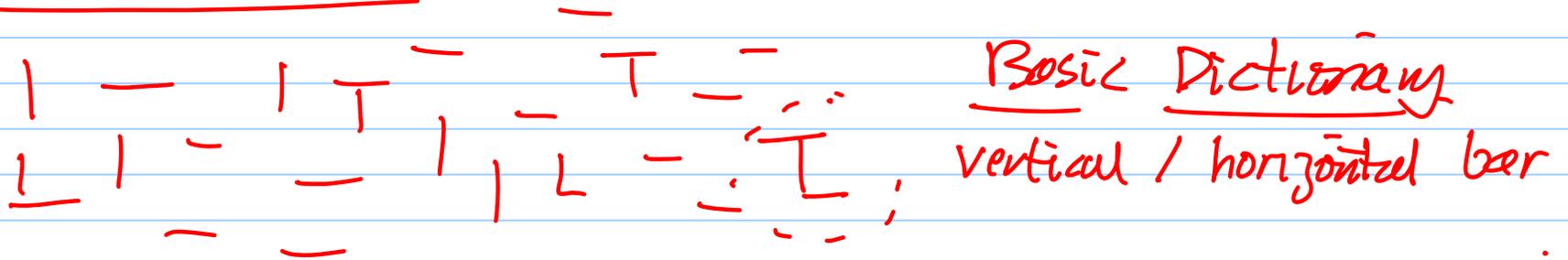
E.g. Cars from multiple viewpoints: Motlaghi & Yuille

Search over Space of Graphical Models: Kemp & Tenenbaum (PNAS)

(13) (II) Unsupervised Learning of Representation & Hierarchies

L. Zhu et al. 2008, 2010

Hierarchies of Edges



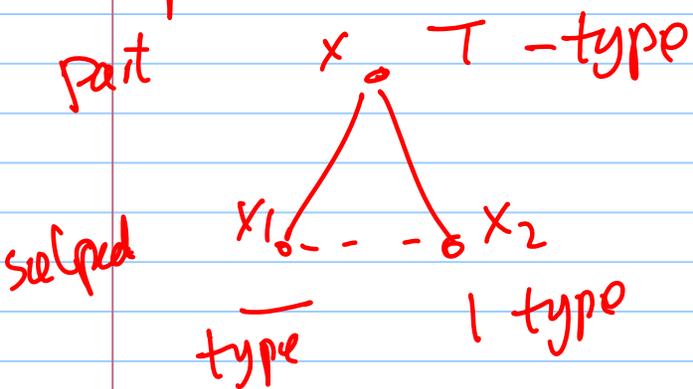
First Stage: same as before

cluster: identify compositors T and L .
don't do model selection \rightarrow don't force decision on T
require $T + B$, $L + B$ to encode better than B

(14) II: Unsupervised Learning of Representations and Hierarchies

Composition: Part-Subpart

S. Geman
Mamford & Desolneux



$$P(x_1, x_2 | x) = \delta(x - \frac{1}{2}(x_1 + x_2))$$

$$h(x_1, x_2; \lambda)$$

$$h(x_1, x_2; \lambda) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x_1 - x_2 - \mu)^2}{2\sigma^2}}$$

$$\lambda = (\mu, \sigma)$$

T = "l" + "-" + "spatial relation."
 $\lambda_1 = (\mu_1, \sigma)$

L = "l" + "-" + "spatial relation"
 $\lambda_2 = (\mu_2, \sigma)$

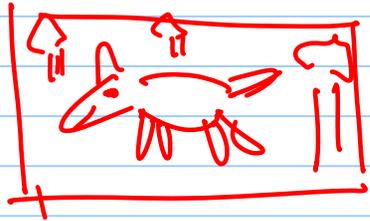
"Visual Concepts"
Deformable Models

(15) (II) : Unsupervised Learning of Repetition and Hierarchies

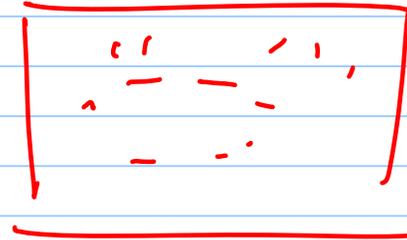
Second Stage:

- Dictionary of low-level features. "—" "I" / , \
- Visual Concepts level 1. — \ T L
- Apply some clustering procedure to visual concepts at level 1 (at this stage resolve ambiguity in T)
- Obtain Visual Concepts level 2.
Repeat → until you stop finding suspicious combinations

(15½) How to learn the representation? (ECCV'08)



→ Edge Detection.
Cocktail Party



Strategy: Low-Level Dictionary | - \ /

Look for compounds of triplets that happen frequently with spatial variability (Gaussian + mean).

• Visual Concepts

• Look for composition of Visual Concepts.

Justification, → Parallel Search over encodings of the image.

Matter of Encoding.

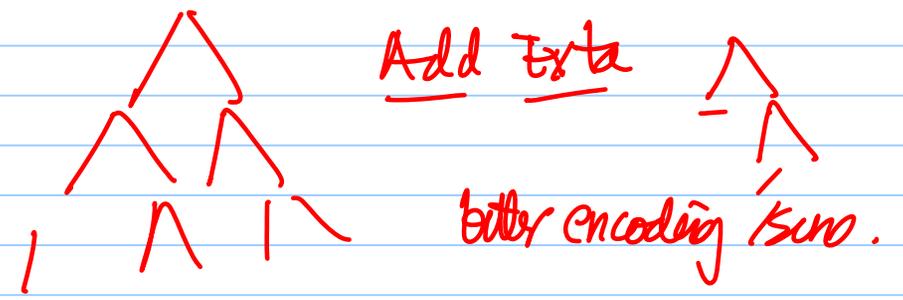
(15 $\frac{3}{4}$) Keep Grouping until you cannot find any components

- Visual Concepts 5
- Visual Concepts 4
- Visual Concepts 3
- Visual Concept 2
- Dictionary 1

Weizman //

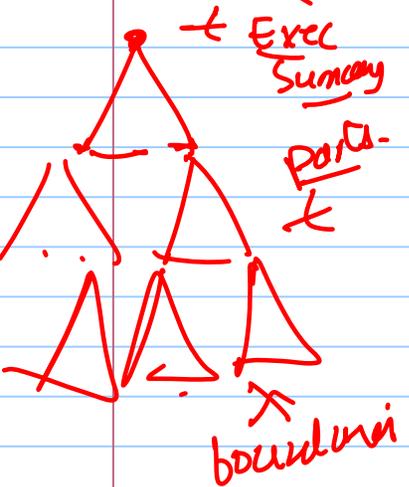
Model. Selection.

Start. with Top Visual Concept

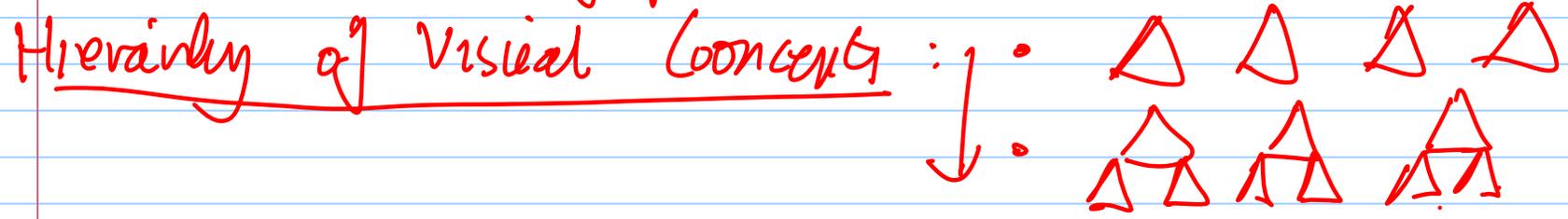


After learning \rightarrow relax $\Delta \rightarrow \Delta$ ^{rotated} _{scale.}
2/3 rule

(16) (II) Unsupervised Learning Representation and Hierarchies



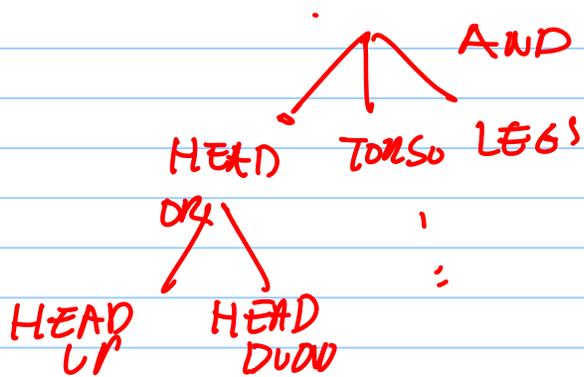
- Learn hierarchical graphical models
- Hierarchical Distributed Representation of objects.
- Coarse - Executive Level Summary
Parent state is invariant to details of children.
- Finer scale detail at lower nodes
positions of parts, boundaries



(17) (II) Unsupervised Learning Representations and Hierarchies

Note: these types of models already existed as hand specified models.

AND-OR Graphs Baseball Players



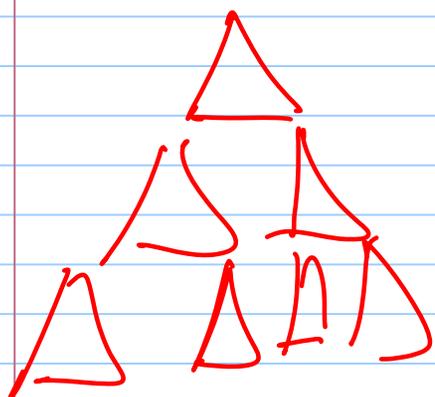
L. Zhu et al.
S. Zhu et al.

S. Geman.
C. Williams.

(+8) (II) Unsupervised Learning of Representations and Hierarchies

Exact Inference on Generative Model.

Feedforward gives executive summary
Feedback resolves low level ambiguities



"Inference on Generative Models can be done very fast."

Feedforward propagates up low-level hypotheses (ambiguous)

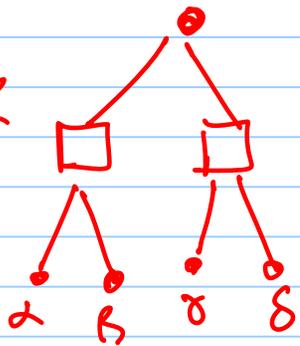
At high-levels there is sufficient context to disambiguate.
Top-Down uses high-level context to resolve low-level ambiguities

(Binding/Linking)

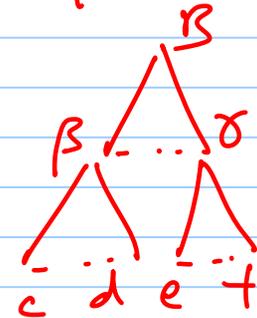
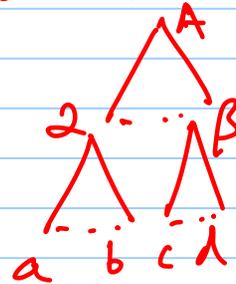
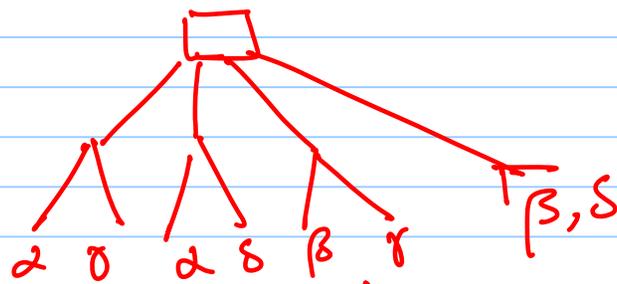
19) (II) Unsupervised Learning Representations and Hierarchies

Convert.

AND-OR graph



→
OR of AND graphs



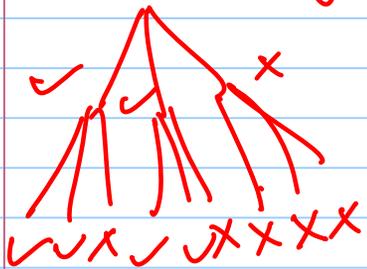
Why Hierarchies?

Part-Sharing.

(20) (III) Complexity Results: Representation and Inference

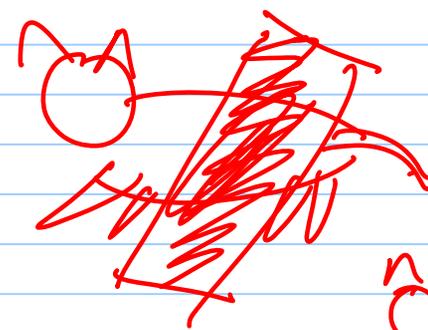
Part Sharing : Efficiency of Representation and Inference (and Learning)

Also robustness to missing parts.



2/3 rule.

only need to detect
4 out of 9 sub-sub-parts



many parts missing at fine level.

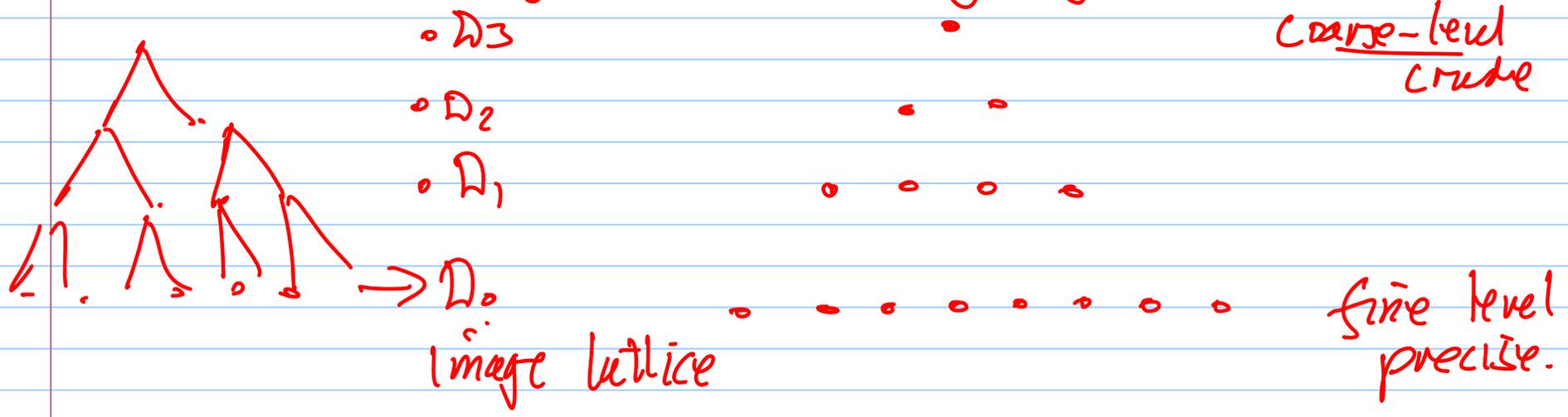
few points missing at high level



(21) (111) Complexity Results: Representation and Inference

Inference: Computation is done by Dynamic Programming
Can compute exact no. of Computations.

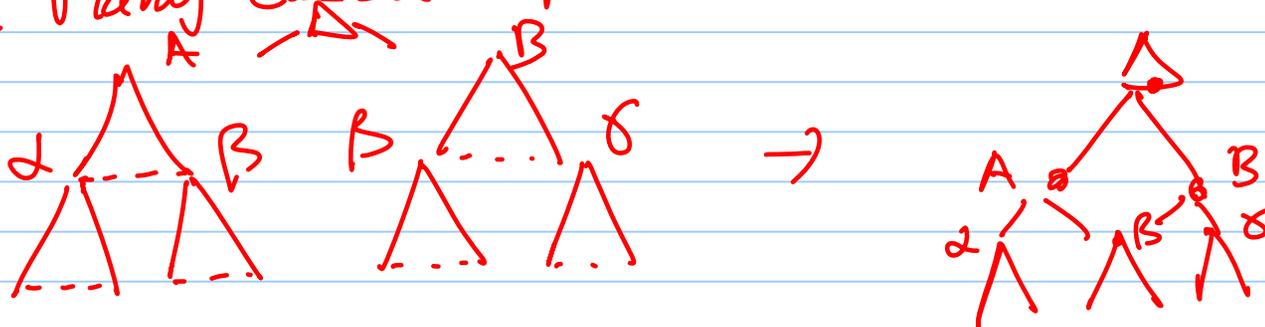
First: "Quantize" the states of object model.



(2 1/2)

Note: Inference is Exact. with Sharing

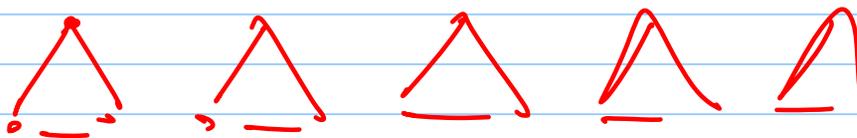
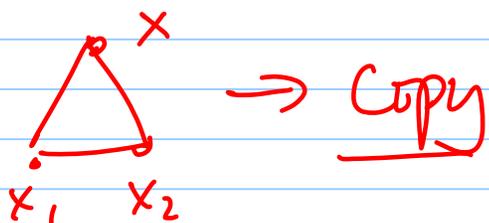
Despite Many Closed Loops.



2. Variant Circuits of the Mind.

Fundamental Problem → Can we derive the structure of the visual cortex from first principles.

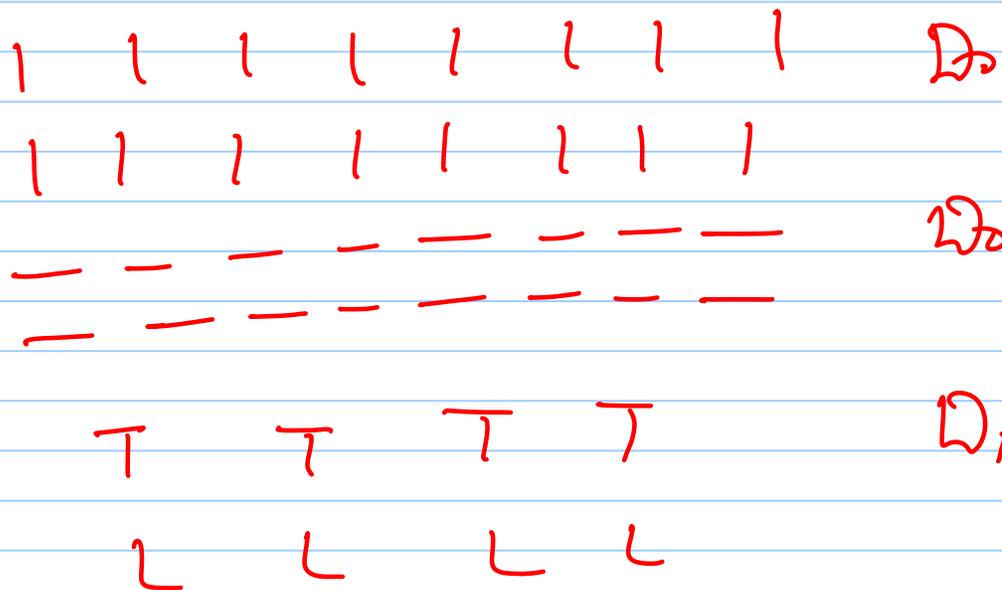
(22) (III) Complexity of Representation and Inference



Lowest-level 1 -

"Columns"

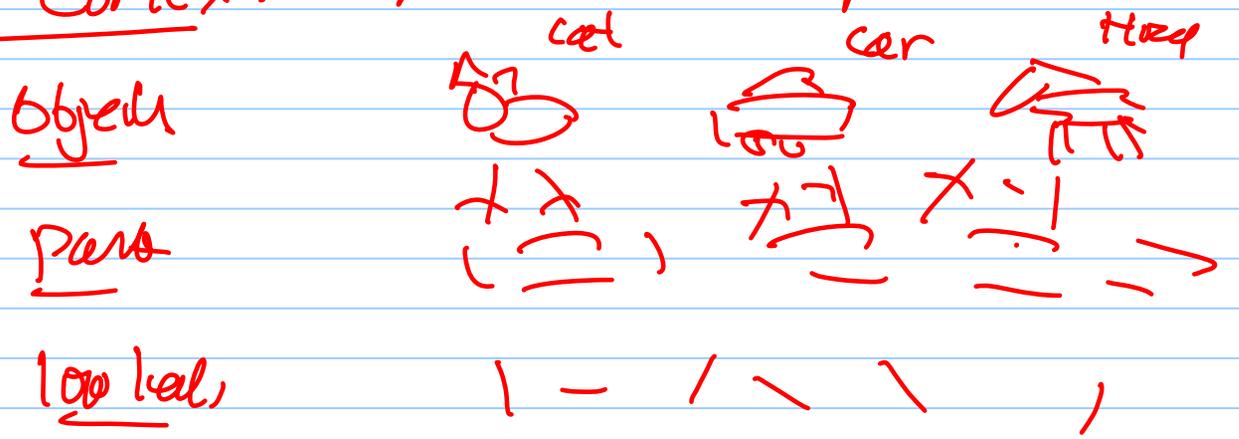
Encode higher-level
concepts sparsely
in position



(23) (III) Complexity of Representation and Inference

Toy Model of Cortex:

Non-linear receptive fields



Sparsely Activated

(24) Complexity: Representation and Inference

Fundamental Problem of Vision → Complexity

How to learn, represent, access all possible objects.
rapidity?

Compositional Hypothesis → possible because objects are
composed in terms of visual concepts (parts) constructed from
more elementary visual concepts (sub-parts).

Hierarchical Representation → Executive Summary
→ Shared Parts.

Visual Architecture → Neural Model.

L. Neuhoff.

• Circuits of the Mind.

(25)

Summary.

- Complexity of Vision: Fundamental Problem?
 - (i) Tasks
 - (ii) Images.
- Compositional Models \rightarrow Explicit Representation of Parts / subparts including spatial relations, Part Sharing
Offers a way to address complexity.
Extensible \Rightarrow Object Appearance, Scenes, 3D, ...

