# Compositional Models: Complexity of Representation and Inference

Beyond Feedforward Models:

"Something about generative models"

(I) General Comments:

Data Driven Models and Representations.

Feedforward and Feedback Architectures

(II) Unsupervised Learning of Representations and Hierarchies.

(III) Complexity Results on Representations and Inference.

(I) General Comments: Data Driven Models.

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

- Face Detection.        (Viola & Jones)
       · Text Detection (Chen & Yuille)

· Pascal Challenge ·  DPM's  (Felzenszwalb, Ramanan, McAllester)
       Object Detection.        · 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.

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

Scene structure: ground plane, Manhattan World structure (d opp?),
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 → 2.5D Sketch → 3D Rep

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

Feedforward and Feedback:  Generative Models

  - Analysis by Synthesis — Mumford & Grenander.

- DDMCMC — Tu & Zhu, Tu, Chen, Yuille Zhu
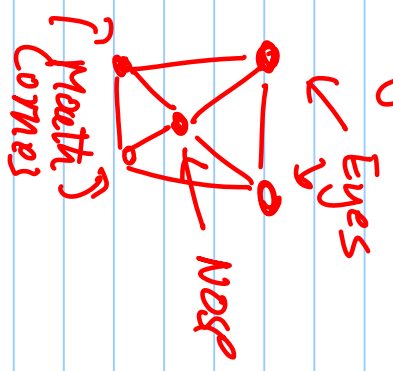  feedforward proposals validated by feedback models.
  But do generative models require feedback ?

  Ullman:
  Streams
  Counter-Streams

(5)

[I] Unsupervised Learning of Representations and Hierarchies

Starting Point : Pictorial Structures. [Fischler & Elschlager]
1973

Model Object (Face) in terms of parts



← Eyes
← Nose
← Mouth) Corners

$$G = (V, \mathcal{E})$$

↗ nodes  ↗ +edges      Graph

State variable $\langle x_\mu : \mu \in V \rangle$

energy  $E[\langle x_\mu \rangle] = \sum_{\mu \in V} \lambda^\mu \cdot \phi(x_\mu) + \sum_{(\mu,\nu) \in \mathcal{E}} \lambda^{\mu\nu} \cdot \psi(x_\mu, x_\nu)$

↗ data      ↗ spatial
term.        relations

Deformable Template.

(6) (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 — Feature Detection "Neural"
              Wiscott, Neven

          • Grylon, Zürich — Hands (Dynamic Programming)

          o Felzenszwalb, Ravmanan, McAllester — Star Model, DPM
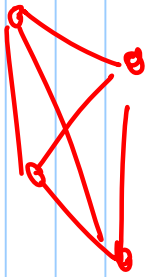              (Pascal Detection)

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

Unsupervised Learning

Constellation Models.          Caltech.

     - Weber , Perona

     • Fergus et al.

Fully connected ?

Flat Models → no composition

# (II) Unsupervised Learning of Representation and Hierarchies

## 1. Zhu et al :

Learn Mixture of Models
Data with clutter

- Take Image → Extract and Represent Interest Point's (IP's)



- Interest Point types. 田 • X
  (~200).

- Task - Learn a Generative Model.

  Don't know : (i) how many objects
  (ii) how many IP's in an object
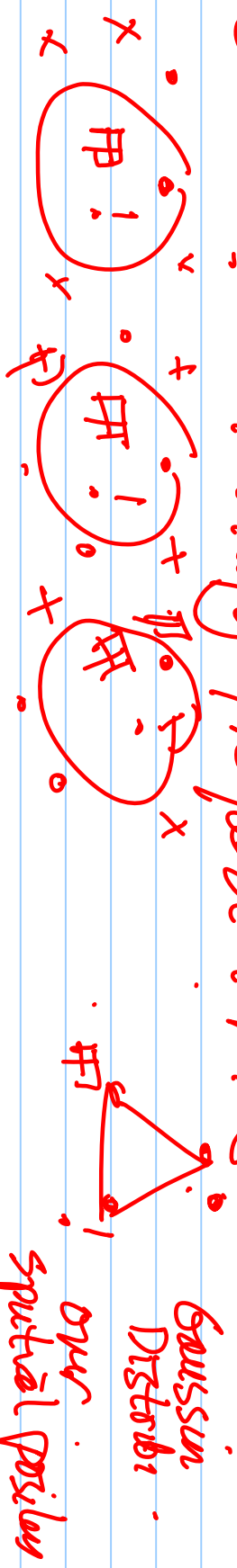  (iii) which IP's are object or background.

# (q) (II) Unsupervised Learning of Representation and Hierarchies

**Strategy:** Greedy search over space of models

**Initialize:** Default model all $IP$'s generated independently $(1,1,1)$

**Cluster:** Identity frequent triplets



Gaussian Distrib.

over spatial posibn

Gaussian over internal angle
invariant to rotation/scale.

**Note:** "Suspicious coincides. — Barlow Model.
" meaningful alignments" Ullman.
" — . . . "

" Suspicious coincides."
" meaningful alignments"
" . . . "

# [9t] (II) Unsupervised Learning of Representation and Hierarchies

## Grow Model by Adding Triplets:

### Current Model:

Image: N background point B → IP's generated by i.i.d background
+ 3M points generated by △ Triplet Model → by or triplet.

is △

Rep by B + △

Grow to B + △ or B + △ + △

↑ more complex object

↗ Grow to more object

# Unsupervised Learning of Representation and Hierarchies

## Better Encoding of Data:

- Cost of Encoding Data by B only: Default

$$\sum \log p(\cdot)$$

- Cost of Encoding Data by $B \rightarrow \Delta$

$$\sum \log p. \ldots + \underbrace{T}_{\text{Overhead}}$$

## Model Selection:

Stop when adding "new object"

or growing b,

(12)

(II) Unsupervised Learning of Representations and Objects

Learnt Representations → Unsupervised
                                       → IP's only  (less interesting)

Learning in Cocktail Party.
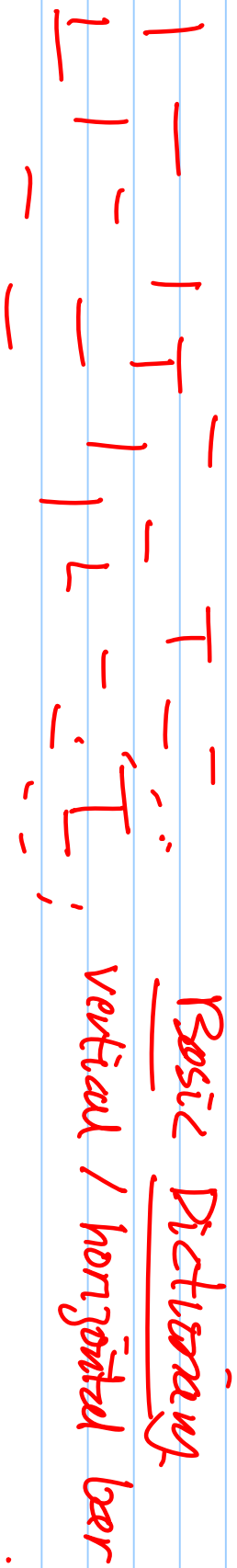       Unknown no. of speakers + Background Noise.

Can do non-trivial Visual tasks

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

       Search over space of Graphical Models :   Kemp & Tenenbaum
                                                    (2008)

(13) (II) Unsupervised Learning of Representation & Hierarchies.
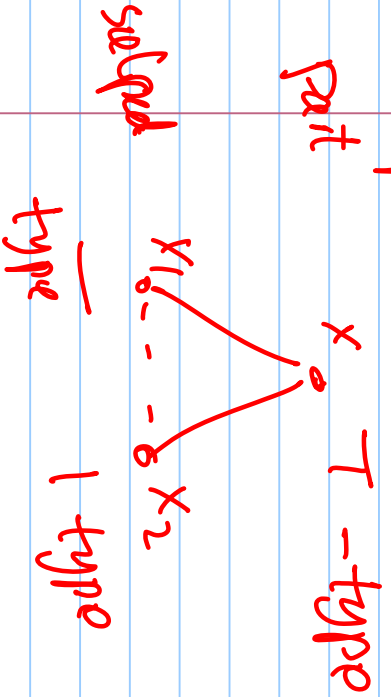
2. Zhu et al. 2008, 2010

Hierarchies of Edges

Basic Dictionary

vertical / horizontal bar

First Stage : same as before

cluster : identity compositons 丁 and L.

don't do model selection → don't force decision on 丁
require 丁 + B , L + B to encode better than B

# (14) (II): Unsupervised Learning of Representations and Hierarchies

**Comparison:** Part-Subpart

Part    x    T-type

        S.Geman

        Mumford & Desolneux

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

subpart $x_1$ ... $x_2$

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

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

1 type

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

$$T = \text{"1" + "—" + "spatial relation"}$$

$$\lambda_1 = (\mu, \sigma)$$

$$L = \text{"1" + "—" + "spatial relation"}$$

        "Visual Concepts"

        Deformable Models

$$\lambda_2 = (\mu, \sigma).$$

(15) (II): Unsupervised learning of Representations and Hierarchies

Second Stage:

- Determining of low-level Features.

- Visual Concepts Level 1.

$$\quad "\_" \quad "\|" \quad / \quad \setminus$$

- Apply Same clustering procedure to visual concepts at level 1 (at this stage resolve ambiguities in $\top$)

- Obtain visual Concepts level 2.
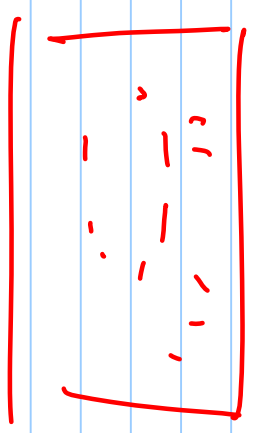
Repeat → until you stop finding suspicious commonalities

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



→ Edge Detection.

Co-lateral Parts

Strategy : Low-Level Dictionary   | ─ ⌣ ✓

Look for complexity of triples that happen frequently  ,  ─
                                          with spatial variability
                                          ( Gaussian + mean ).

• Visual Concepts

• Look for complexity of visual Concepts .

Justification, → Parallel Search over encodings of the image .      Matter of
                                                                    Encoding.

(15 3/4) <u>Keep Grouping until you cannot find any component</u>.

<u>Visual Concepts</u> 5

<u>Visual Concepts</u> 4

<u>Visual Concept</u> 3

<u>Visual Concept</u> 2

Pictures 1

Neurmin //

<u>Model Scheduler</u>.

<u>Start with</u> Top Visual Concept

<u>Add text</u>

Better encoding func.

<u>After learning</u> → $\dfrac{relax}{2/3 \ rule}$

$\triangle \rightarrow \triangle$ \ rotate scale.

(II) Unsupervised Learning Representation and Hierarchies

+ Exec Summary → Learn hierarchical Graphical models

• Hierarchical Distributed Representation of object.

- Coarse → Executive Level Summary .·. Parent state is <u>invariant</u> to details of children.

- Finer scale details at lower nodes positions of parts, boundaries

Hierarchy of visual Concepts

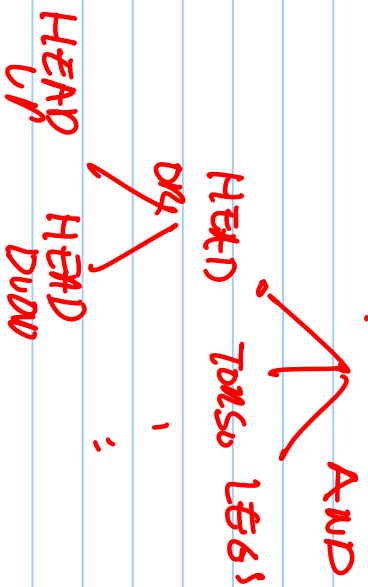(17) (II) Unsupervised Learning: Representations and Hierarchies

Note: these types of models already existed
as hand specified models.
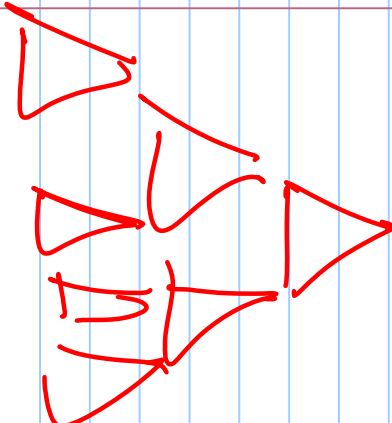
AND-OR Graphs: Baseball Player

L. Zhu et al.

S. Zhu et al.

S. Geman.

C. Williams.

HEAD TORSO LEGS

AND

HEAD
Up

HEAD
Down

(+2) (II) Unsupervised Learning of Representations and Hierarchies

Exact Inference on Generative Model.

Feedforward gives executive summary/
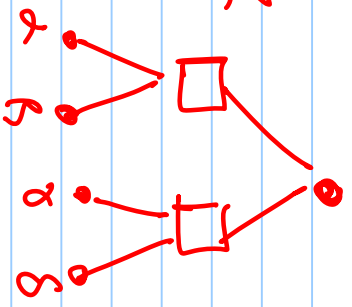Feedback resolves (too) level ambiguities

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

Feedforward propagates up low-level hypotheses (ambiguous)
At high-levels that is sufficient context to disambiguate.
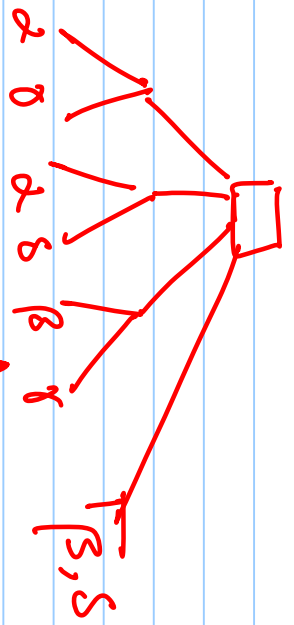Top-Down uses high-level context to resolve low-level ambiguities
(Binding / Linking)

# (II) Unsupervised Learning Representation and Hierarchies
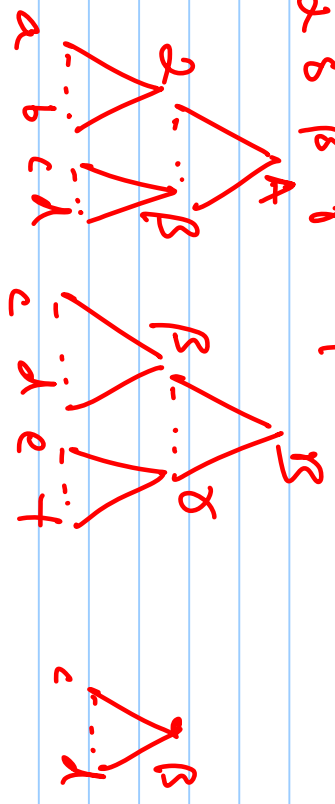
**Convert.**
AND-OR Graph  →  OR of AND Graphs
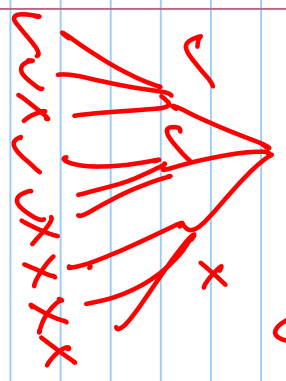
**Why Hierarchies?**

Part-Sharing.

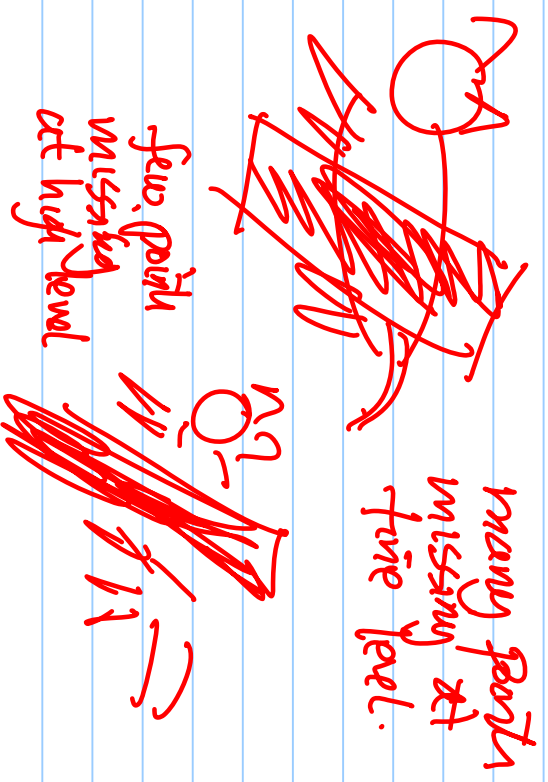[20] (II) Complexity Results: Representation and Inference

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

Also robustness                                          many parts
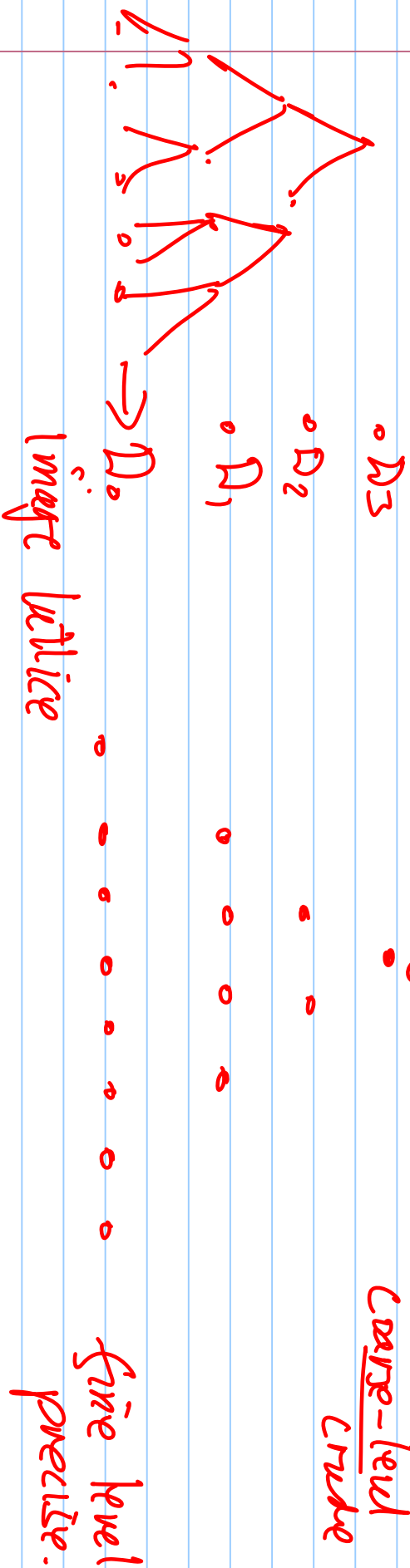to missing parts.                                        missing at
                                                         true level.

+      2/3 rule.

only need to detect                     few parts
4 out of 9 sub-sub-parts                missing
                                        at high level.

(21) (iii)

## Complexity Results: Representation and Inference

Inference:

First: "Quantize" the states of object model.

Computation is chosen by Dynamic Programming
can compute exact no. of computation.

• D3          coarse-level
                crude
• D2

• D1

→ D0
      fine level
      precise.

Image lattice

Note : Influence in Exact · with Shapes

Despite Many Closed Loops.



→
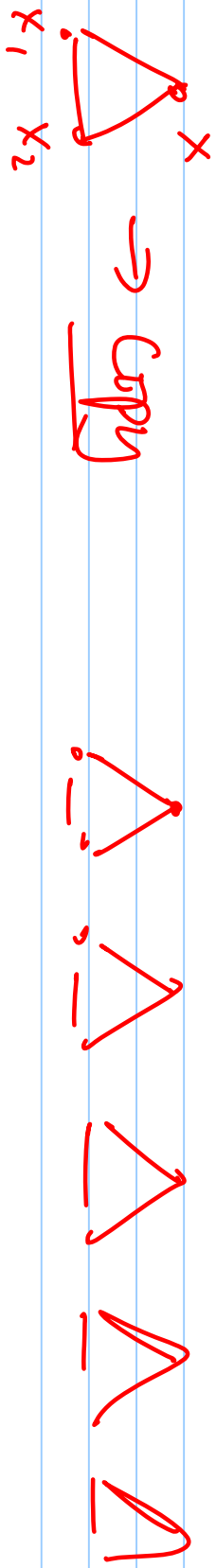
2. Valiant — Circuits of the Mind.
Fundamental Problem → Can we derive the structure of
the visual cortex from first principles.

## (II) Complexity of Representation and Inference

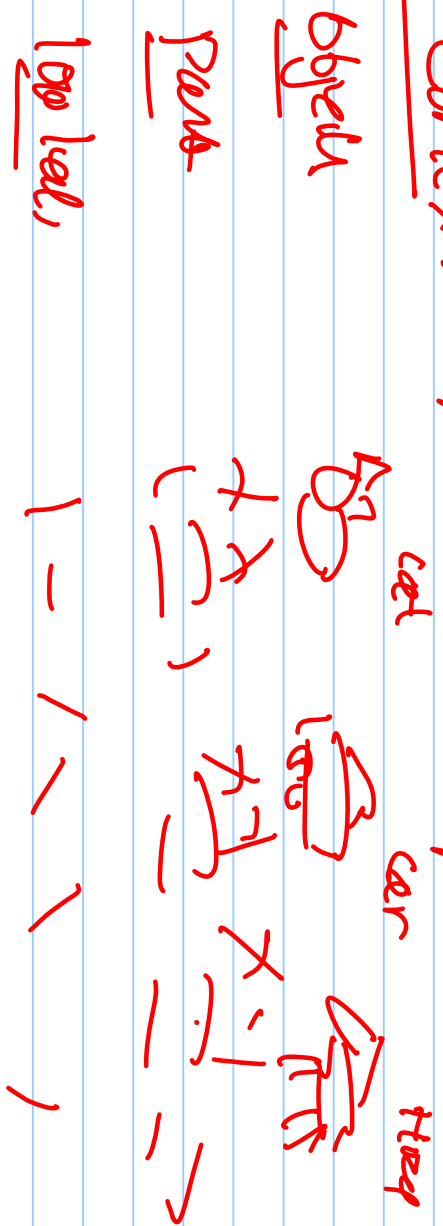$\triangle^{x}_{x_1 \ x_2} \rightarrow$ Copy

Lower-level )

"Columns"

Encode higher-level

Concepts Sparsely
in position

$D_2$

$D_0$

$D_1$

## (23) (II) Complexity of Representation and Inference

**Toy Model of Context:** Non-linear receptive fields

objects

Pens

low level

(24)

# Complexity: Representation and Inference

## Fundamental Problem of Vision → Complexity.

How to learn, represent, access all possible objects rapidly?

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

Hierarchical Representation → Executive Summary
→ Shared Parts.

Visual Architecture → Neural Model.
1. Voltaic.
   Circuit of the Mind.

(25)

# Summary.

• Complexity of Vision:       Fundamental Problem?

 (i) Tasks

 (ii) Images.

• Compositional Models → Explicit Representation of Part/Subparts including spatial relations; Part Sharing.

   offers a way to address complexity.

       Extensible ⇒ Object Appearance, Scenes,
                      3D, .