

OPTICAL CHARACTER RECOGNITION

Motivating Example: OCR

Optical Character Recognition:

- image → text
- closely related to object recognition, autonomous driving, speech recognition, ...
- much better statistical understanding of...
 - class structure
 - noise
 - ground truth
 - syntactic structure / language modeling

Applications

- surprisingly, a hot topic!
- wanted: information extraction from unstructured documents
- NLP technology didn't use to be up to it
- with new LLMs, that has changed dramatically

Motivating Example: OCR

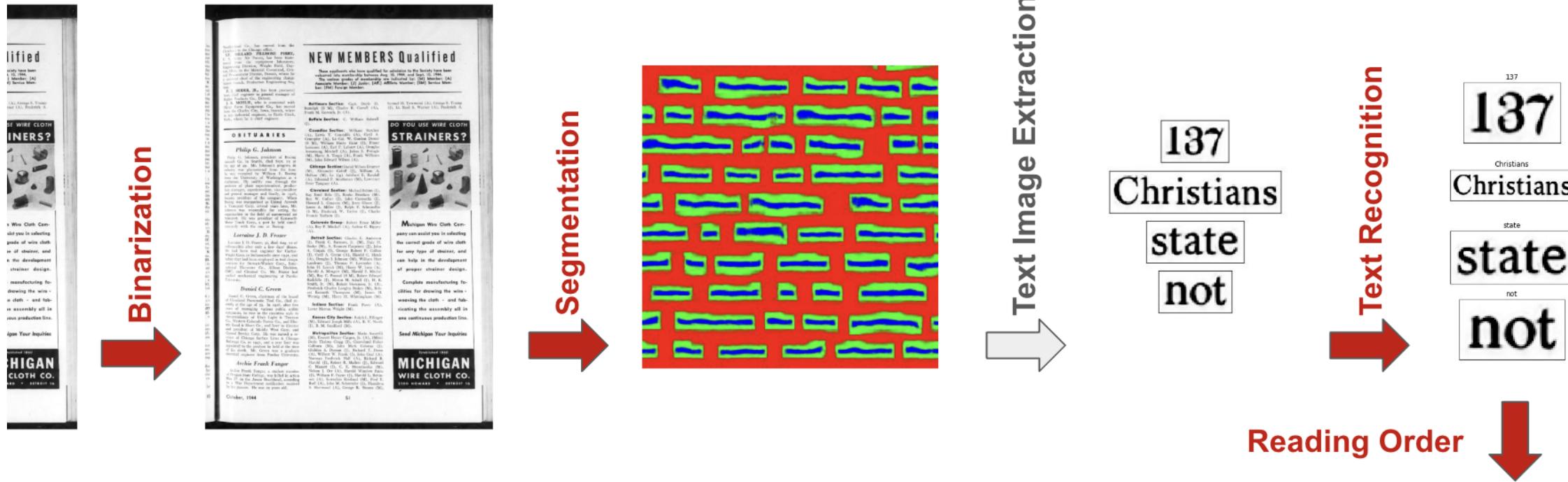


A horizontal right-pointing arrow indicating the direction of the next section.

```
<html version="1.0" encoding="UTF-8">
<!DOCTYPE html PUBLIC "-//W3C//DTD XHTML 1.0 Transitional//EN"
<!http://www.w3.org/TR/xhtml1/DTD/xhtml1-transitional.dtd">
<html xmlns="http://www.w3.org/1999/xhtml" xml:lang="en" lang="en">
<head>
<title></title>
<meta http-equiv="Content-Type" content="text/html; charset=utf-8"/>
<meta name="ocr-system" content="tesseract 4.1.1-r25-g9707" />
<meta name="ocr-capabilities" content="ocr_page ocr_carea ocr_par ocr_line ocr_word ocrp_wconf"/>
</head>
<body>
<div classe='ocr_page' id='page_1' title='image "sample.jpg"; bbox 0 0 2130 3433; ppageno 0'>
<div classe='ocr_carea' id='block_1_1' title='bbox 115 224 1884 284'>
<span classe='ocr_par' id='par_1_1' lang='eng' title='bbox 115 224 1884 284'>
<span classe='ocr_line' id='line_1_1' title='bbox 115 224 1884 284; baseline -0.001 -11; x_size 51.400451; x_descenders 6.400456; x_ascenders 18'>
<span classe='ocr_word' id='word_1_1_1' title='bbox 115 224 228 284; x_wconf7 1654>'>
<span classe='ocr_word' id='word_1_1_2' title='bbox 598 228 824 284; x_wconf7 93'>Prerogative</span>
<span classe='ocr_word' id='word_1_1_3' title='bbox 867 223 979 275; x_wconf7 71'>Court</span>
<span classe='ocr_word' id='word_1_1_4' title='bbox 1822 232 1884 271; x_wconf7 96'>of</span>
<span classe='ocr_word' id='word_1_1_5' title='bbox 1583 231 1845 271; x_wconf7 73'>Canterbury.</span>
<span classe='ocr_word' id='word_1_1_6' title='bbox 1738 236 1884 271; x_wconf7 87'>+273</span>
</span>
</span>
</div>
</div>
<div classe='ocr_carea' id='block_1_2' title='bbox 118 292 1792 381'>
<span classe='ocr_par' id='par_1_2' lang='eng' title='bbox 118 292 1792 381'>
<span classe='ocr_line' id='line_1_2' title='bbox 118 292 1792 381; baseline 0 3132; x_size 28; x_descenders 5'>
<span classe='ocr_word' id='word_1_2_1' title='bbox 118 292 1792 381; x_wconf95'> </span>
</span>
</span>
</div>
<div classe='ocr_carea' id='block_1_3' title='bbox 116 358 1885 399'>
<span classe='ocr_par' id='par_1_3' lang='eng' title='bbox 116 358 1885 399'>
<span classe='ocr_line' id='line_1_3' title='bbox 116 358 1885 399; baseline -0.003 -12; x_size 41; x_descenders 8; x_ascenders 12'>
<span classe='ocr_word' id='word_1_3_8' title='bbox 116 352 376 389; x_wconf90'>SLADDEN,</span>
<span classe='ocr_word' id='word_1_3_9' title='bbox 402 353 540 395; x_wconf94'>James,</span>
<span classe='ocr_word' id='word_1_3_10' title='bbox 566 352 604 395; x_wconf93'>of</span>
<span classe='ocr_word' id='word_1_3_11' title='bbox 626 354 666 395; x_wconf91'>19</span>
<span classe='ocr_word' id='word_1_3_12' title='bbox 689 353 727 395; x_wconf96'>St.,</span>
<span classe='ocr_word' id='word_1_3_13' title='bbox 770 353 808 397; x_wconf96'>Margaret,</span>
<span classe='ocr_word' id='word_1_3_14' title='bbox 808 354 846 395; x_wconf96'>Rochester,</span>
<span classe='ocr_word' id='word_1_3_15' title='bbox 1118 353 1339 395; x_wconf96'>Kent,</span>
<span classe='ocr_word' id='word_1_3_16' title='bbox 1364 351 1478 398; x_wconf96'>July</span>.
<span classe='ocr_word' id='word_1_3_17' title='bbox 1560 354 1698 395; x_wconf96'>July</span>.
<span classe='ocr_word' id='word_1_3_18' title='bbox 1612 360 1750 395; x_wconf96'>27,</span>
<span classe='ocr_word' id='word_1_3_19' title='bbox 1766 350 1885 392; x_wconf96'>+1653.</span>
</span>
</span>
```

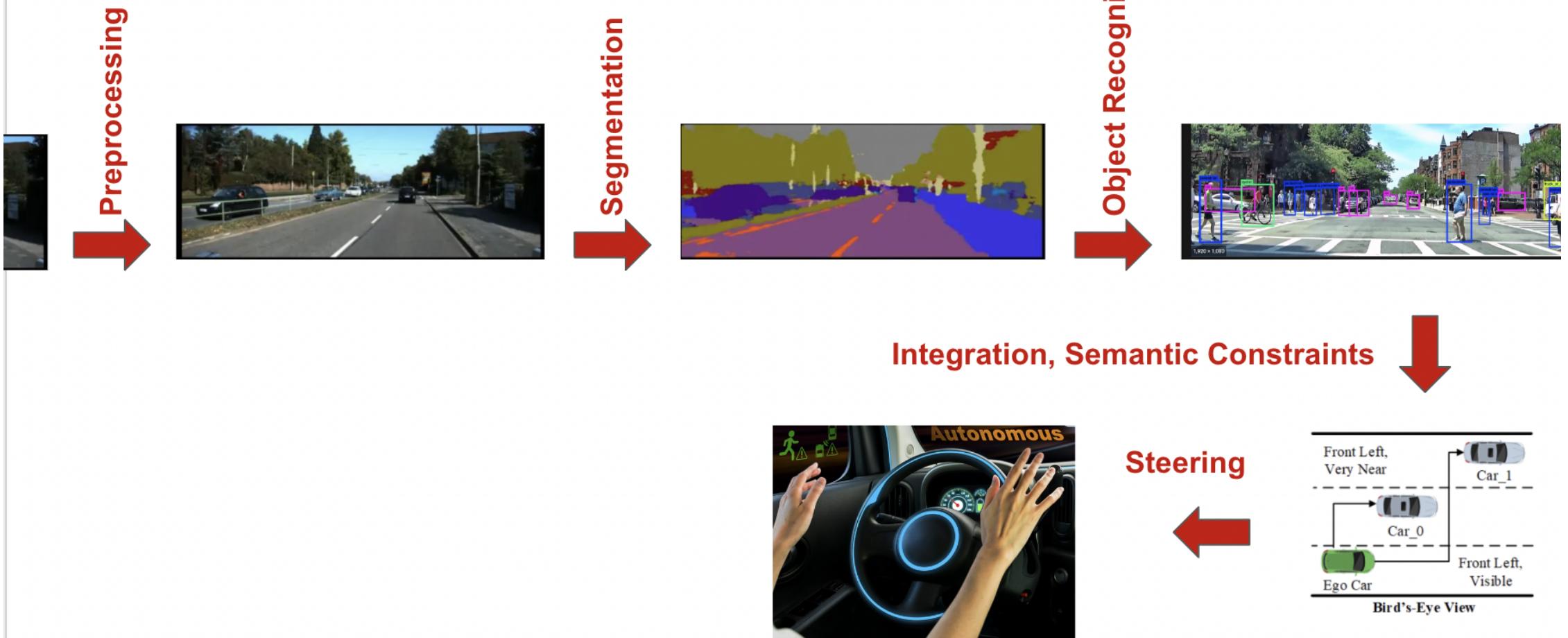
image → text, markup

Motivating Example: OCR



1500s, when an unknown printer took a galley of type and scrambled it to make a type specimen book. It has survived not in their exact original form, reproduced in their exact original form,

OCR vs AV Pipeline



Motivating Example: OCR

We need to train multiple models:

- preprocessing
- segmentation
- recognition
- syntactic model (language model)
- geometric model (reading order model)

This is common to many computer vision tasks using deep learning.

OCR and Semi-Supervised Learning

Available data:

- thousands of scanned pages with manual segmentation and text
- tens of thousands of scanned pages with approximate text
- millions of scanned pages without text
- large amounts of text without document images
- the ability to generate new texts and printed documents

Typical unsupervised / semi-supervised learning problem.

First Step: Supervised Training

Train supervised models on manually labeled / transcribed data:

- segmentation model
- image → text model

Yields good performance *on data similar to training data.*

Unlabeled training data primarily helps with generalization to new datasets.

Using Millions of Untranscribed Pages

Idea:

- Use our or other OCR systems to transcribe those pages and use the output as training data ("pseudolabel")

Questions:

- How can using a worse OCR system work for training a better OCR system?
- Is there anything we can do to improve this?

SEMI-SUPERVISED TEXT RECOGNITION

Example: Word Recognition Problem

Let's focus just on recognizing words for the following examples (forget about the rest of the OCR system):



137
Christians
state
not
what

Using Millions of Untranscribed Pages

Idea:

- run the existing OCR system, giving word images and corresponding text
- construct a new training set by...
 - rejecting any word that the OCR system has a low confidence in
 - rejecting any word that is not found in the dictionary
- this way, we obtain a new training set that contains "mostly good" training samples
- iterate this multiple times

Why does this work?

Network estimates $P(c|x)$ (x : image, c : class)

For OCR, we know:

- true $P(c|x)$ is approximately 0 or 1 (no ambiguities)
- any uncertainty in classifier output is due to mislabeled training data
- e.g. if 20% of training data mislabeled: $\tilde{P}(c|x) = 0.8$ for true class c
- if we use pseudolabel $\arg \max_c \tilde{P}(c|x)$, many accidentally mislabeled training labels will actually be correctly labeled
- as a result, the posterior probability will be estimated higher on the next training round and the model improves

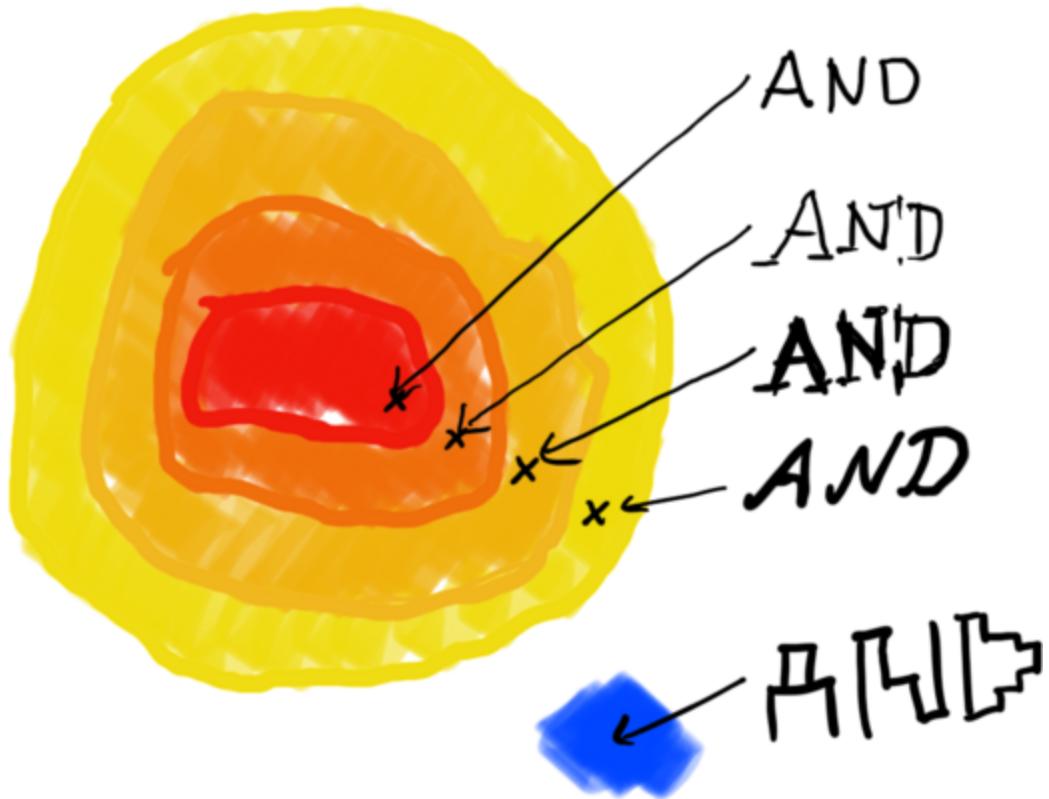
EM Algorithms

- latent variables in semi-supervised OCR
 - labels for unlabeled portion of training set
 - outlier status for unlabeled portion of training set
- EM algorithms recover latent variables by...
 - "making a best guess" given the current model
 - retraining the model as if that guess is correct

Unsupervised Learning

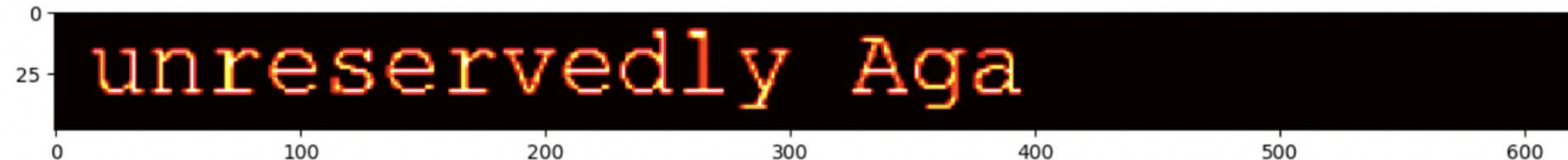
- identify the *latent variables* that are being recovered
- identify the *prior assumptions / inductive biases* in the model
- identify the *EM algorithm implementation*

Iterated Recognition / Dataset Construction

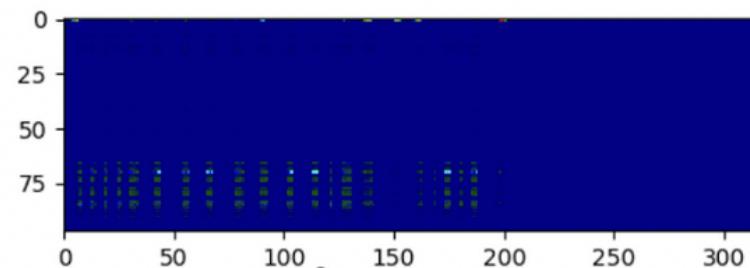


WEAKLY SUPERVISED TEXT RECOGNITION

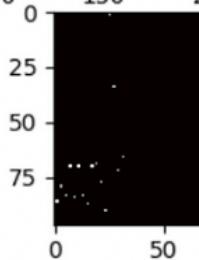
EM Training with Language Models (aka CTC)



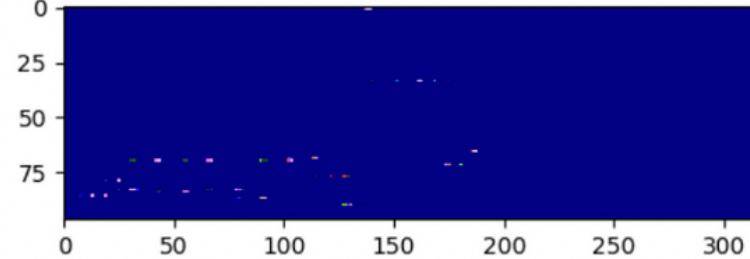
posterior probability at each pixel location



one-hot ground truth



CTC-aligned ground truth



EM Training with LSTM/Conv + CTC

- ideally, the LSTM/Conv model outputs the correct class only where the character occurs
- outputs ϵ (no character) everywhere else
- our transcript does not contain the character positions
- CTC performs *alignment* between classifier output and transcript
- CTC estimates the most likely positions of characters given the current model



EM Training without Any Transcripts

Normal EM-Training:

- image + transcript, EM-training recovers alignment

Language Model-Based EM-Training:

- image + language model, EM-training recovers text + alignment

EM Training, Information Theory

There always has to be *some* source of information:

- unigram perplexity: 1000 ← what the classifier outputs
- bigram perplexity: 200
- trigram perplexity: 100 ← what the language model imposes
- full transcript perplexity: 1 ← fully supervised training

Additionally: several bits per character for character location.

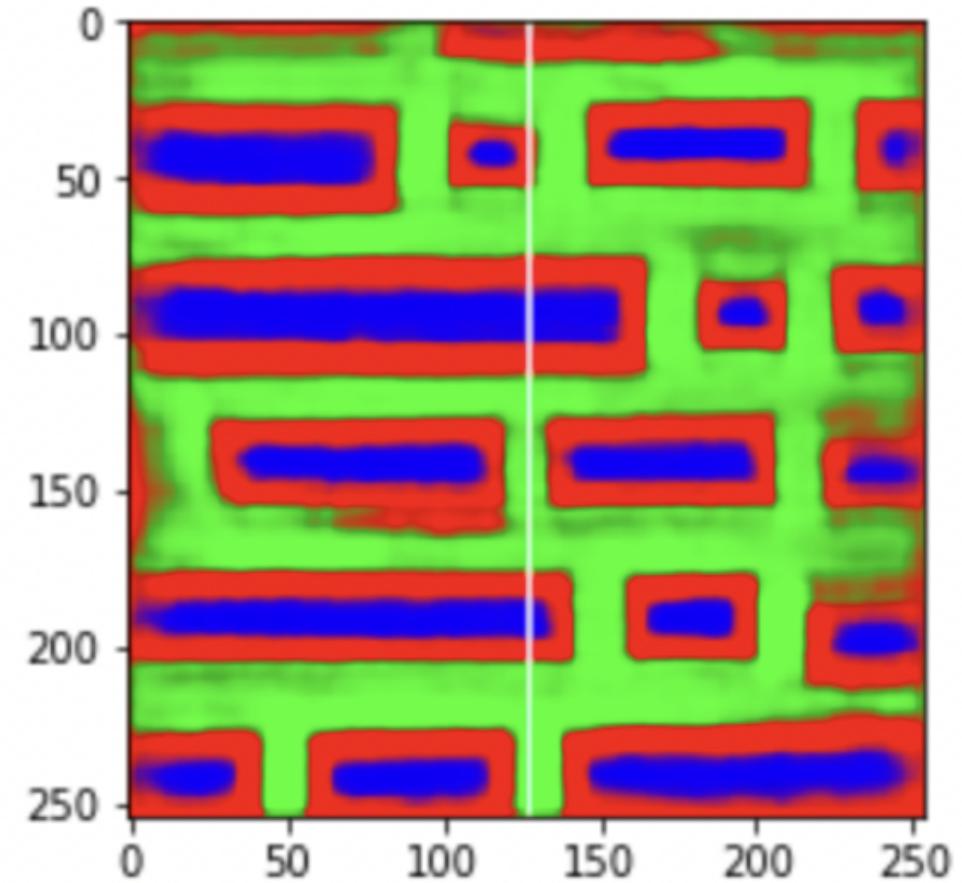
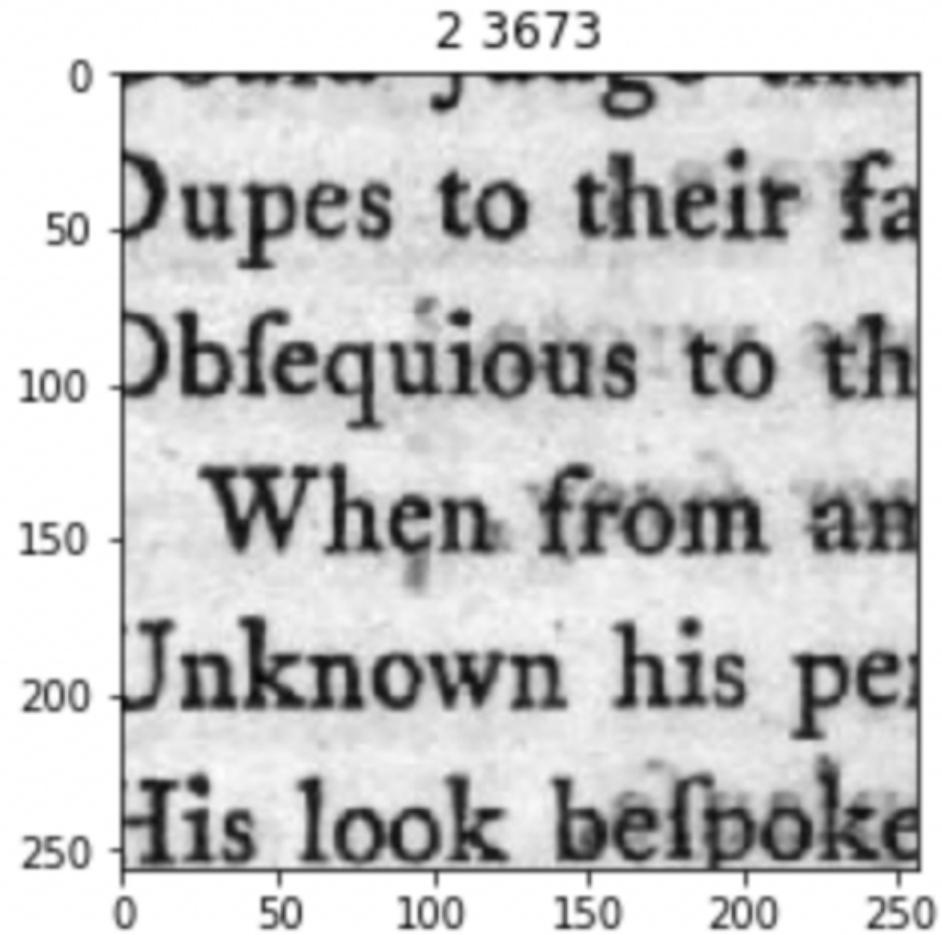
Different Forms of Weak-/Unsupervised Learning

We have seen three different forms of weak/unsupervised training for text → image training:

- using image + transcript, lacking just the alignment/segmentation
- using image + language model only
- using just pseudolabels from weak classifier

SELF-SUPERVISED PAGE SEGMENTATION

Self-Supervision of Page Segmentation



Self-Supervision of Page Segmentation

Assume we start with errorful segmenter from small labeled training data.

Two kinds of errors:

- segmenter returns word for non-words
- segmenter misses actual words

Approach:

- validate each returned word using OCR (no OCR = not a word)
- mark everything that is not identified as a word as "don't know" and exclude from training

Self-Supervision

Approaches to self-supervision usually incorporate prior knowledge and require some manual design:

- text recognizer
 - reject non-words from "soft labels"
 - use language models as part of EM training
- page segmentation
 - validate segmentation via OCR
 - introduce "don't know" mask during training

Self-Supervision

Prior knowledge about the task is required to choose meaningful self-supervision tasks.

- word recognition: pseudolabels, language-model based rejection
- segmentation: verification via OCR, introduction of "don't care" regions
- object recognition, natural image segmentation: *later*

ACTIVE LEARNING

Active Learning

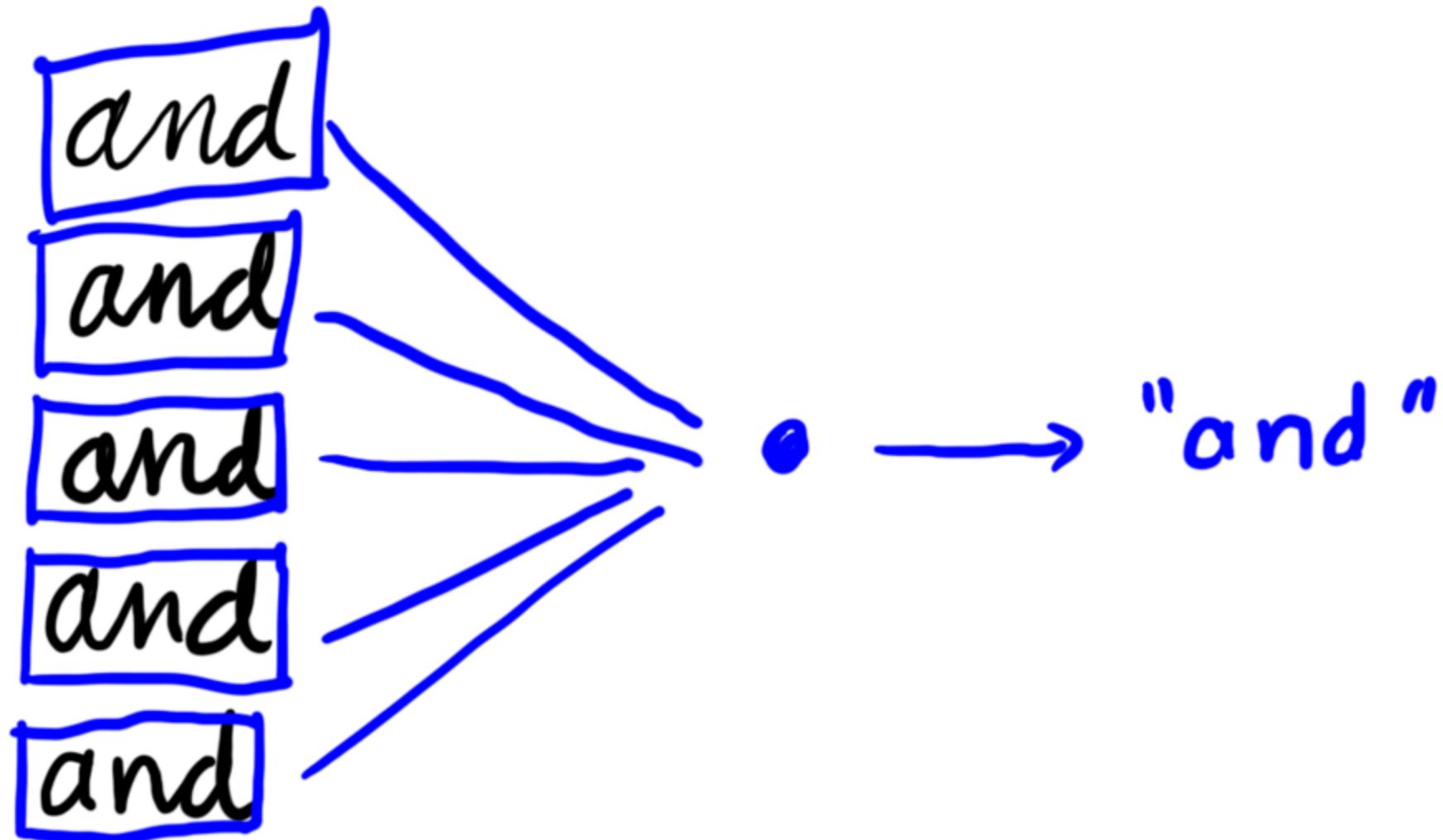
Request Help from Oracle:

- run OCR as above and send low confidence outputs to manual transcribers
 - output near decision boundaries
 - output that can maximally help the classifier improve

Clustering:

- perform clustering on character or word images and manually transcribe each cluster

Transcription with Clustering



Question: what similarity measure do we use for clustering?

OCR by Solving a Cryptogram

ယူစွဲ သတ္တမ မိန္ဒီ ဒြပ် နဲ့ ဖျောက် နှင့် လုပ်ချက် နဲ့ ယူစွဲ မိန္ဒီ ပြုပေး
နဲ့ လျော် မိန္ဒီ ဒြပ် ယူစွဲ မိန္ဒီ ပြုပေး လုပ်ချက် နဲ့ ပြုပေး

- cluster characters by shape
- use frequency and patterns to infer character identity
 - e.g.: first word is likely "THE"
 - letter frequencies
 - word frequencies
 - single letter word is likely "A"
- explicit form of many EM-based recognition algorithms

OCR by Solving a Cryptogram

THE CAT SAT ON A WALL AND LOOKED AT THE SUN.
A DOG SAT ON THE STREET AND CHEWED A BONE.

DATA SOURCE MODELS

Hierarchical Bayes



- clustering is based on a simple hierarchical Bayesian source model
- latent variables can be recovered by EM or Bayesian methods
- simple clustering assumes $\Sigma = 1$
- both c and μ are latent; need oracle to recover actual class labels

Channel View

$$c \in \mathbb{Z}_k \xrightarrow{\text{render}} \xi \xrightarrow{\text{noisy channel}} x \xrightarrow{\text{classifier}} P(c|x)$$

OCR: render = digital typeset; noisy channel: printing + scanning

Vision: render = choice of pose ; noisy channel: lighting, camera

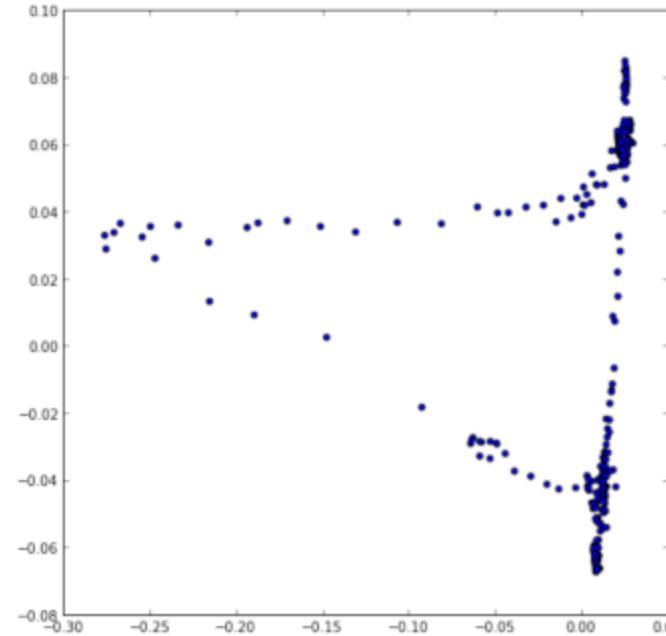
Rendering can involve transformation parameters, giving rise to *view manifolds*.

Real View Manifold



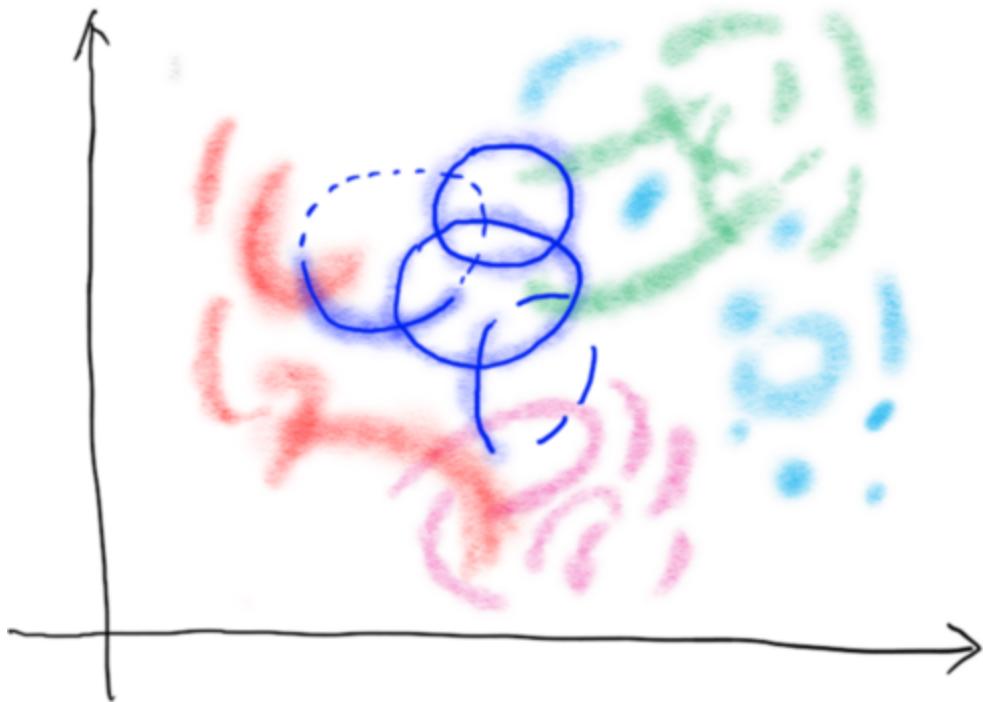
5 second motion sequence of object rotating
back and forth

```
from sklearn.decomposition import RandomizedPCA
pca = RandomizedPCA(20)
lo = pca.fit_transform(data)
mds = manifold.LocallyLinearEmbedding()
vl = mds.fit_transform(lo)
scatter(vl[:,0],vl[:,1])
```



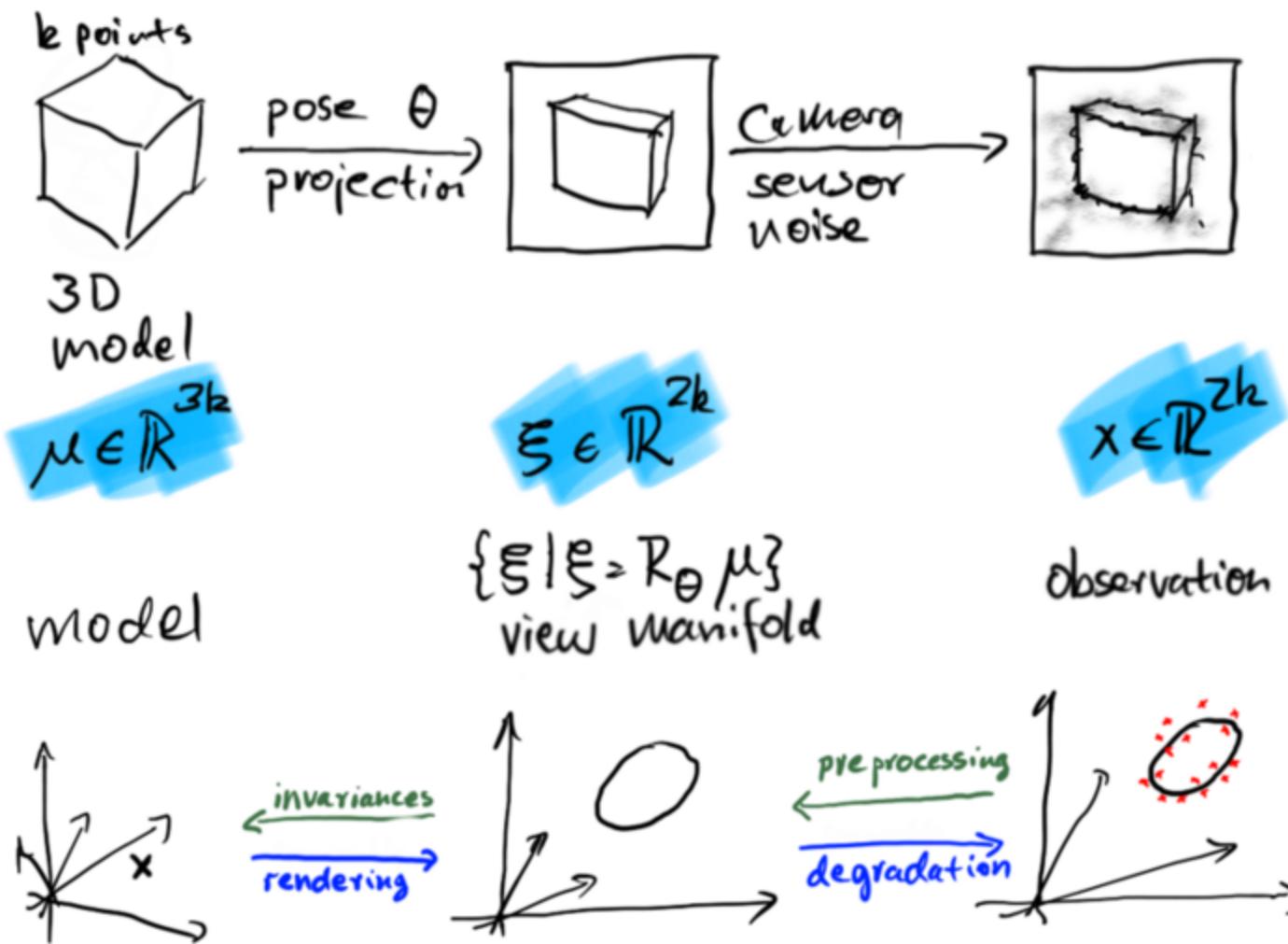
Dimensionality reduction in pixel space: $\mathbb{R}^{256} \rightarrow \mathbb{R}^2$. Video: [rotating.mp4](#)

Manifolds Learning



- think of real data as clouds surrounding...
- union of manifolds with boundary (only shown for blue class)
- "manifold learning" = union-of-manifold-s-with-boundary learning

Model, View Manifold, and Noise



Unsupervised Learning

Given or Learned:

- rendering
- degradation

Constructed or Learned (inverse problems):

- invariances
- preprocessing / image cleanup

CycleGAN can learn both directions simultaneously with no supervision.

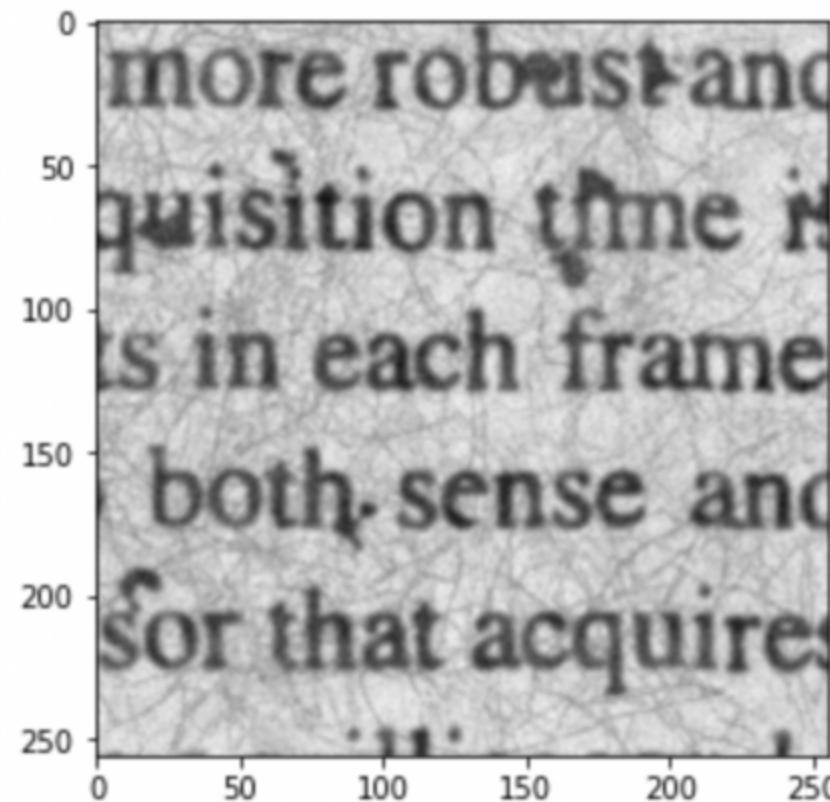
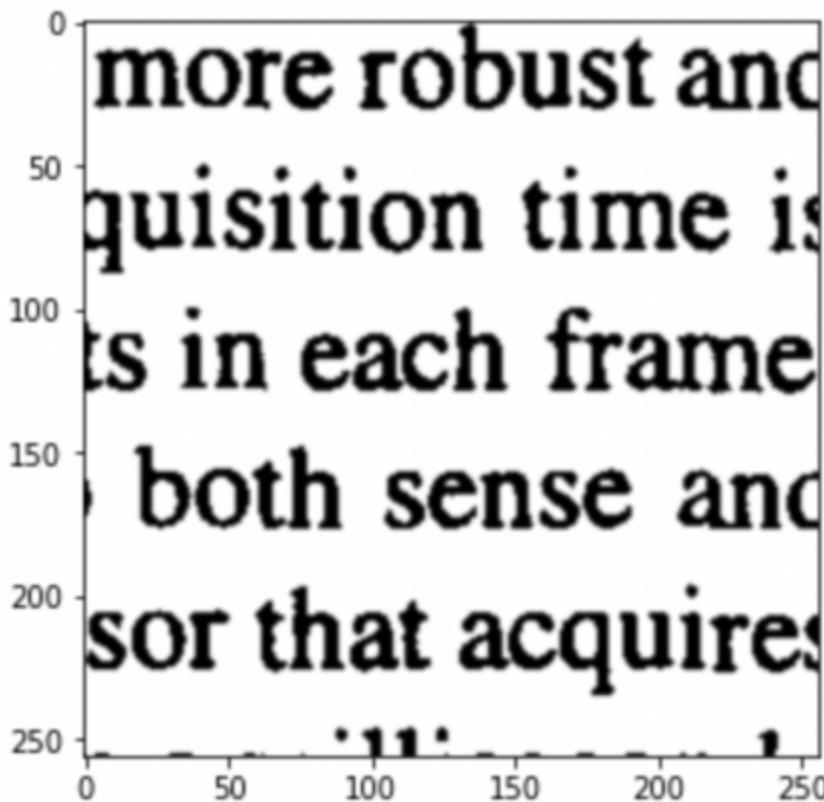
TRAINING DATA GENERATION / DEGRADATION

Artificial Data Generation

This is the "forward path" in the channel view of recognition:

- digital typesetting = perfect "artificial" document generation
- take any text, generate pages of perfect text
- for OCR, we need "degraded images"
 - printing, scanning, photographing, photocopying, ...
- physical processes for document image degradataion are well known

Artificial Data Generation



Document Image Restoration

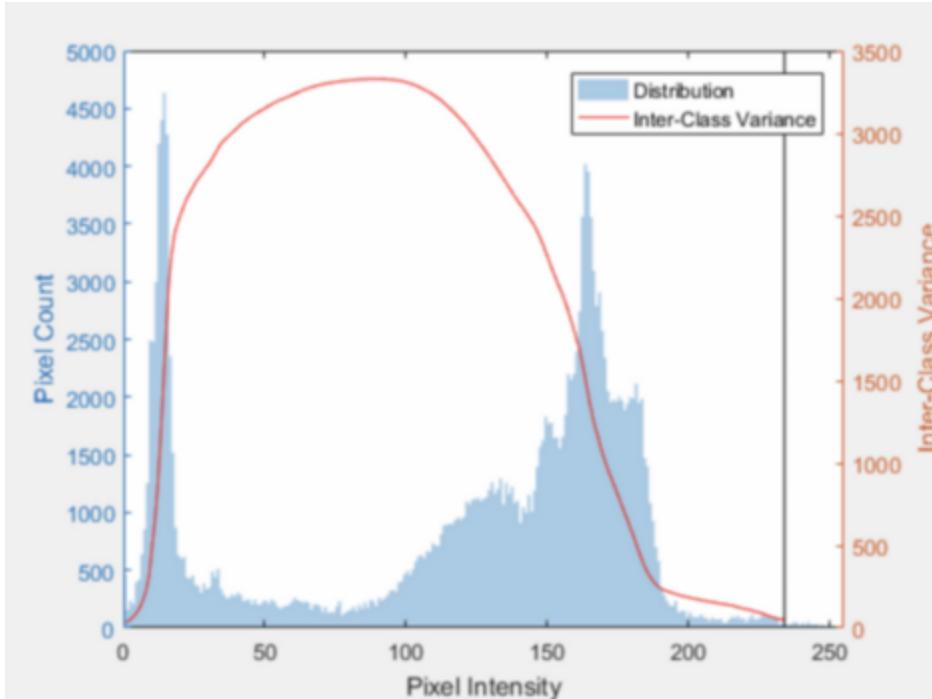
- recognition is easier if documents are not degraded
- can we restore the "clean" image via unsupervised learning?

AV/object recognition: image translation prior to recognition

Classical Document Image Restoration

- binarization
 - optimal thresholding
 - optimal linear filtering (deconvolution, etc.)
 - clustering
 - performance-based dynamic thresholding
- deep learning
 - supervised restoration (LSTM, pix2pix)
 - unsupervised restoration via CycleGAN

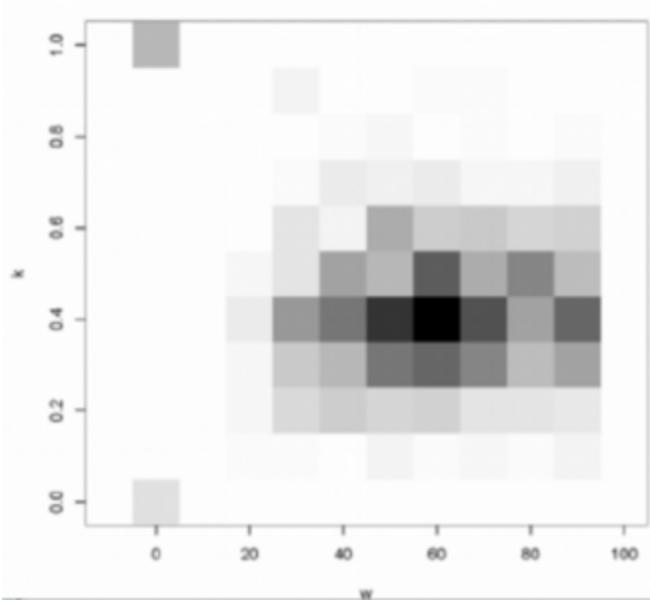
Otsu's Method: "Optimal" Thresholding



Otsu's method:

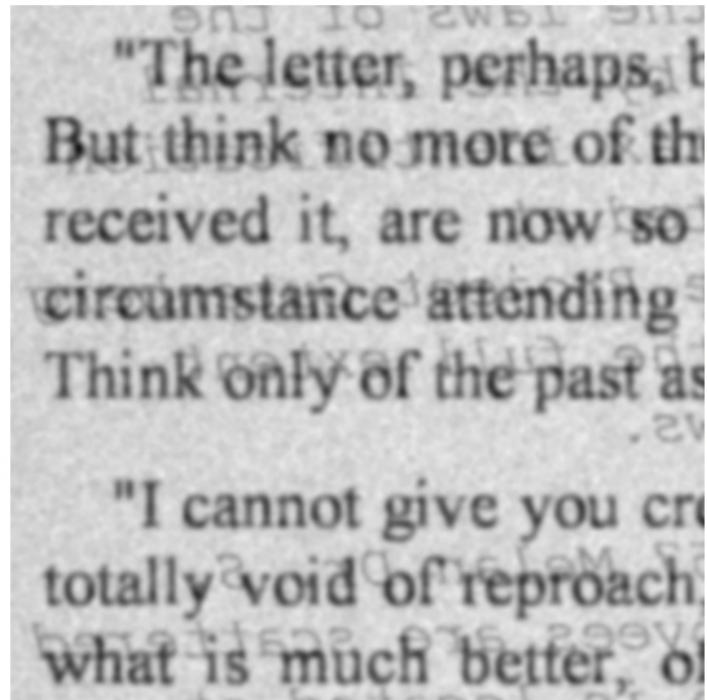
- manually constructed inverse degradation model based on prior knowledge/assumptions
- assume some kind of mixture model of image pixel generation
- maximize inter-class variance

Performance-Based Thresholding



- run thresholding with many parameters
- pick the parameters that yield the best OCR output
- no transcript, so use proxy
 - use statistics/classifier, or #words in dictionary
 - closely related to GAN methods
- assumption/prior: local thresholding = good way of inverting degradation model

Binarization by Supervised Deep Learning



(a)

"The letter, perhaps, it
But think no more of the
received it, are now so
circumstance attending
Think only of the past as

"I cannot give you cre
totally void of reproach,
what is much better, ol

(b)

"The letter, perhaps, it
But think no more of the
received it, are now so
circumstance attending
Think only of the past as

"I cannot give you cre
totally void of reproach,
what is much better, ol

(c)

(a) original, (b) Sauvola, (c) LSTM-based binarization

Self-Supervised Training for Binarization

1. generate clean images, degrade with degradation model
 - takes advantage of prior knowledge of degradataion models
2. take degraded images, use performance-based thresholding
 - takes advantage of knowledge of statistical properties of output

CycleGAN for Document Preprocessing

CycleGAN replaces all those components with trainable networks:

- clean image → degraded image
- degraded image → clean image
- clean image detector
- degraded image detector

These correspond to the forward and backward arrows in our channel model.

CycleGAN can be trained end-to-end without any labeled data or (significant) prior assumptions.

CycleGAN for Document Preprocessing

EMOIR OF THE AUTHOR.

his thought, that nothing is i
; of Christ to engage in, i
effectually promote the ki
laker. Perhaps it is not p
d the world will hardly belie
een taken in composing th
what care I have endeav

 d_A

$$g_{AB} \rightarrow$$

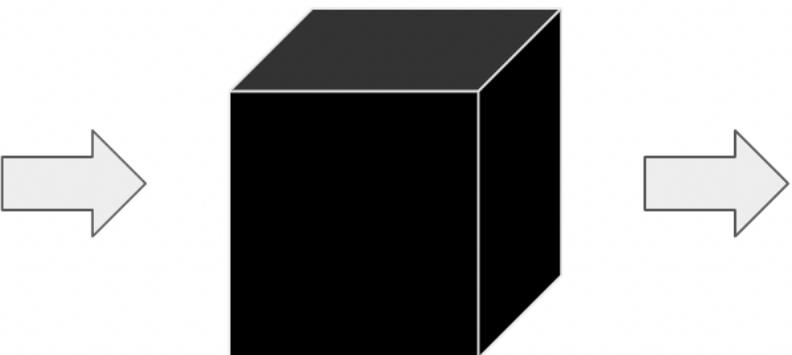
$$\leftarrow g_{BA}$$

EMOIR OF THE AUTHOR.

his thought, that nothing is i
; of Christ to engage in, i
effectually promote the ki
laker. Perhaps it is not p
d the world will hardly belie
een taken in composing th
what care I have endeav

 d_B

End-to-End OCR



Weatherhead Co., has moved from the Cleveland to the Chicago office.

LT. MILLARD FILLMORE PERRY, U. S. Army Air Forces, has been transferred from the equipment laboratory, Engineering Division, Wright Field, Dayson, Ohio, to the Materiel Command, Central Procurement Division, Detroit, where he is assistant chief of the engineering change branch.

...

End-to-End OCR

Recent Developments

- transformer models permit full end-to-end training for image-to-text transcriptions
- do not need intermediate segmentation, text-line recognition

LayTR Multicolumn Recognition

Image-to-HTML tags.

LayTR Table Recognition

tweezed	hydroxylase	preinvasive	filices	unlabialising
annalists	implodes	lachrymation	soir	scleroprotein
yashiro	zibeth	vespering	crotchets	ascescency
hundredponny	trigonoid	aquatical	limburgito	prododication
similarize	rhizocephala	brooklike	supervalued	contrary
milligramage	cativo	yorkshireism	protension	geleem
midnightly	pedatifid	scienced	bacteriosis	germanomania

*fleahopper gallivanting
oncoming heathenize
breakback celebrators
overwrestle sampler*

bayberries	cyanidin	nicotinise	stercorary	incogitant	talepyet
unlethal	snuffle	poltroonery	scimiterpod	abasedness	boomers
sediments	mulctatory	spindlewise	timocratical	renunciatory	standardizer

midianitish | elastomeric
rodding | polymelian
cottonwick | verbena
smouser | gingalls
shippound | verifiably
melter | lurch

Figure 3: Four samples (a)-(d) featured in the LayTr-Tables dataset.

Image-to-HTML tags.

Self-Supervised Training

- apply the same principles
- train an initial model using supervised data
- compute output on unlabeled data ("soft labels")
- correct the output using language models
- retrain
- possibly use auxiliary tasks for pretraining (later)

Or has it been solved already...

Summary

In the OCR example, we have seen most of the major concepts of unsupervised and semi-supervised training:

- EM training (aka soft labeling)
- use of language modeling as data source
- clustering
- unsupervised preprocessing / image enhancement