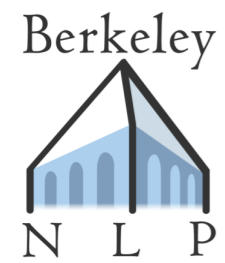# Natural Language Processing



## Historical Document Transcription

Dan Klein — UC Berkeley

Joint work with Taylor Berg-Kirkpatrick and Greg Durrett [ACL 2013]

# Historical Document

## Old Bailey Court Proceedings 1775

the prisoner at the bar. Jacob Lazarus and his wife, the prisoner, were both together when I received them. I sold eleven pair of them for three guineas, and delivered the remainder back to the prisoner. I sold seven pair of silk to Mark Simper: one pair of mixed, and two pair of thread to the footman, and one pair of thread to the barber.

Q. What is the footman's name?

*Frances Moses.* I don't know.

*Henry Harris.* I was standing at the Compter waiting for the sherriff's officers to employ me : Moses's daughter came for me to go and take the prisoner. I went to the Old Bailey

# Transcription

## Document Image

the prisoner at the bar. Jacob Lazarus and his wife, the prisoner, were both together when I received them. I sold eleven pair of them for three guineas, and delivered the remainder back to the prisoner. I sold seven pair of silk to Mark Simper: one pair of mixed, and two pair of thread to the footman, and one pair of thread to the barber.

2. What is the footman's name?

*Frances Moses.* I don't know.

*Henry Harris.* I was standing at the Compter waiting for the sherriff's officers to employ me : Moses's daughter came for me to go and take the prisoner. I went to the Old Bailey

Berkeley NLP

## Document Image



## Transcription (Google Tesseract)

modern

m   o

modern

m o d

modern

# Historical Document

the prisoner at the bar. Jacob Lazarus and his wife, the prisoner, were both together when I received them. I sold eleven pair of them for three guineas, and delivered the remainder back to the prisoner. I sold seven pair of silk to Mark Simper: one pair of mixed, and two pair of thread to the footman, and one pair of thread to the barber.

Q. What is the footman's name?

*Frances Moses.* I don't know.

*Henry Harris.* I was standing at the Compter waiting for the sherriff's officers to employ me : Moses's daughter came for me to go and take the prisoner. I went to the Old Bailey

long s glyph

(a)

(b) the Death of the Deceased,

(c)

(a)

(b) the Death of the Deceased,

(c)

Wandering Baseline

(a)

(b) the Death of the Deceased,

(c)

(a)

(b)

(c) rude along in silence

(a)

(b)

(c) rule along in silence

(a)

(b)

(c) rule along in silence

Uneven Inking

# Various Historical Documents

(a) Old Bailey, 1725:

> (b) Old Bailey, 1875: the Prisoner came drunk to his Stand, (at Mr. Bird's Door in Castle-Court) and without any Provocation began to be very quarrelsome, swearing, calling him ill Names, and striking him two or three times; Hemmit desired him to get out of his Beat, or he'd make him forfeit Sixpence. (Such a Forfeit being customary among the Watchmen, if one comes into the other's Beat.) Mr. Bird then came to the Door, and threaten'd the Prisoner that he would charge a Constable with him, and send him to Bridewell; upon which the Prisoner was very free of his ill Language to Mr. Bird,

(b) Old Bailey, 1875:

(c) Trove, 1823:

> in the Police Office at Sydney, on the 7th of November, 1821. Mr. Norton, the plaintiff's solicitor, laid the case before the Judge in nearly the following words.— He stated that his client, in this case, was Mr. James the owner of the schooner Little Mary, a resident at Port Dalrymple, and that the defendant was Mr. Peter Dillon, commander and owner of the late East India ship Fatisalam, with which vessel he sailed from Bengal bound to these Colonies, with a valuable cargo, but

(d) Trove, 1883:

(b) Old Bailey, 1875:

(c) Trove, 1823:

> to be conscious at the time—he was caught up and thrown out between them—he could not resist—that was the last I saw—I then went down stairs, and saw him lying on the stones below the window, and his mother came up directly after—I saw a soldier catch hold of William Bagley—all the men went down stairs directly they had thrown him out of the window —there is a court-way leading from the lodging house into the street—I cannot say what way they went, they walked away quickly—I did not see Mr. or Mrs. Rowe or Joseph do anything to promote this attack.
> *Cross-examined.* There were six or seven men by the window at the time he was thrown out—I was standing by the fire-place—I saw him actually thrown out by the four—I could not see whether he went out head first or
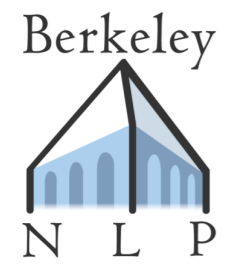
(c) Trove, 1823:
(d) Trove, 1883:

> I have been told. In this case I suspect our voices startled them. We will wrap it carefully as it is."
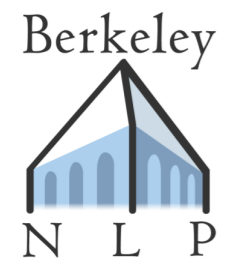> "What for?" said Boolger.
> "The murderers might be iden
> "Hardly. One tomahawk is just
> "Yes, but there are marks on that which, though meaningless to us, are

(d) Trove, 1883:

positive

# Our Approach

(a) *positive*

(b) the Death of the Deceased,

(c)

Berkeley
N L P

(a) positive

(b) the Death of the Deceased,

(c) rule along in silences

# Generative Model

p

```
p r
```

p r i s o n e r

# Generative Model

prisoner

Language Model

```
p r i s o n e r
```

Typesetting
Model

p r i s o n e r

Typesetting
Model

p r i s o n e r

Typesetting
Model

p r i s o n e r

Typesetting
Model

p r i s o n e r

Typesetting Model

p r i s o n e r

Typesetting
Model

p r i s o n e r

Typesetting
Model

p r i s o n e r

Typesetting
Model

# Generative Model

p r i s o n e r

Typesetting
Model

# Generative Model

p r i s o n e r

Typesetting
Model

Berkeley
NLP

p r i s o n e r

Typesetting
Model

Typesetting Model

# Generative Model

p r i s o n e r



Typesetting
Model

p r i s o n e r



Rendering Model

# Generative Model



prisoner

Rendering Model

prisoner



Rendering Model

p r i s o n e r

# Generative Model

Language Model
$P(E)$

$E$

p r i s o n

# Generative Model

Language Model
$P(E)$

$E$

Typesetting Model
$P(T|E)$

$T$

p r i s o n

# Generative Model

Language Model
$P(E)$

Typesetting Model
$P(T|E)$

Rendering Model
$P(X|E,T)$

$E$

$T$

$X$

p r i s o n

# Generative Model

Language Model
$P(E)$

Typesetting Model
$P(T|E)$

Rendering Model
$P(X|E,T)$

$E$

$T$

$X$

$$\boxed{E}$$

Kneser-Ney smoothed character 6-gram

a

$e_i$

a

$e_i$

$T$

Left pad width

Right pad width

Glyph box width
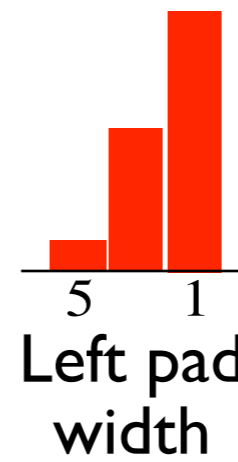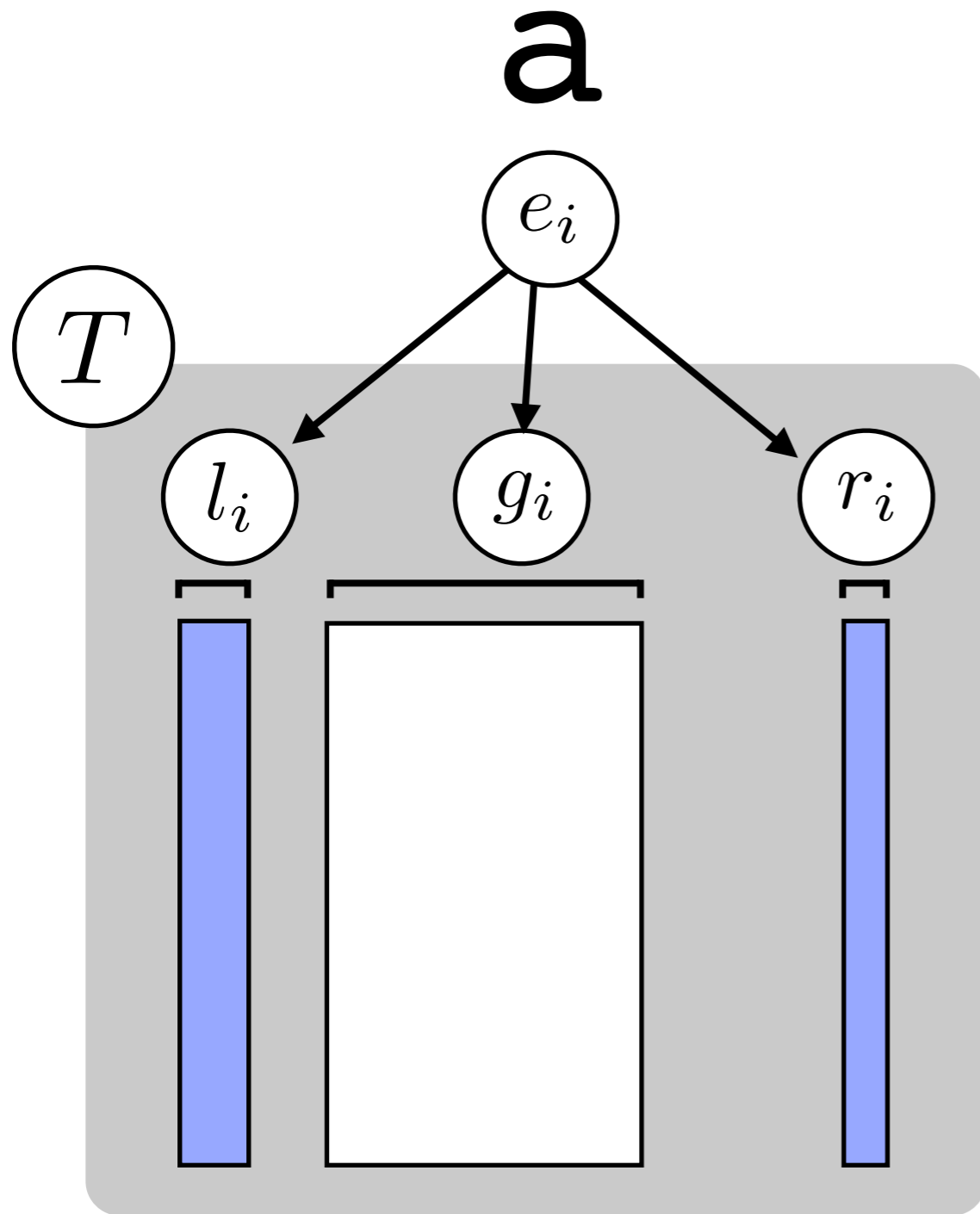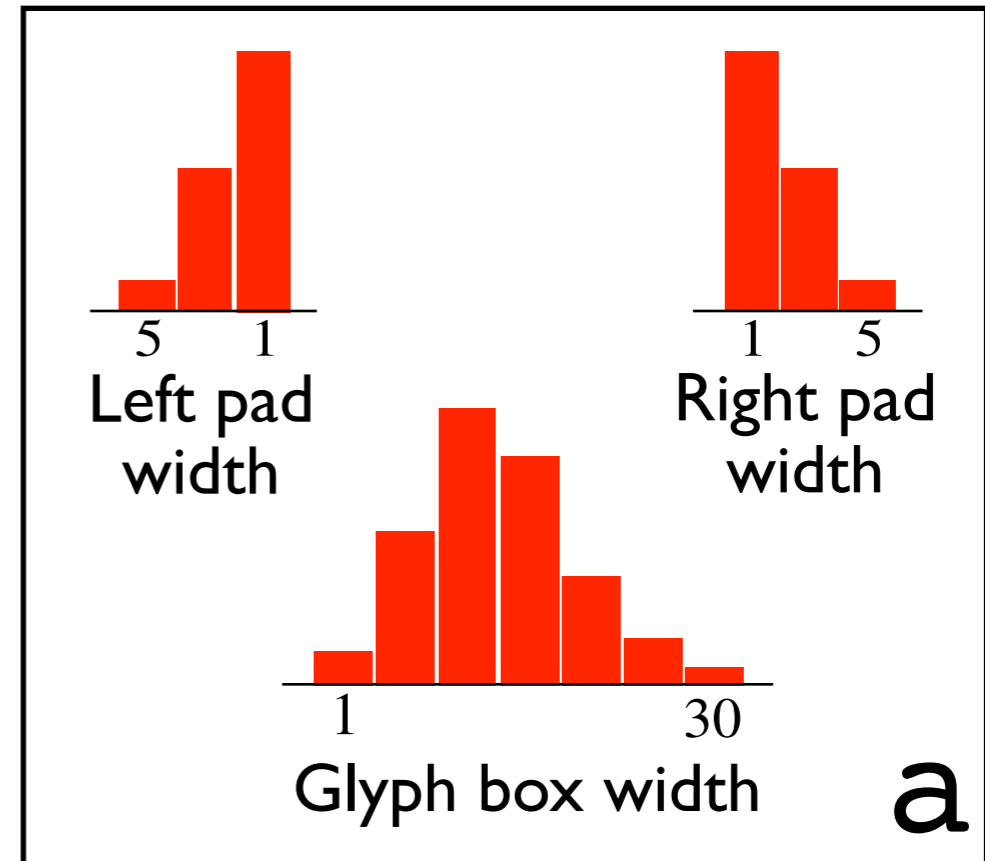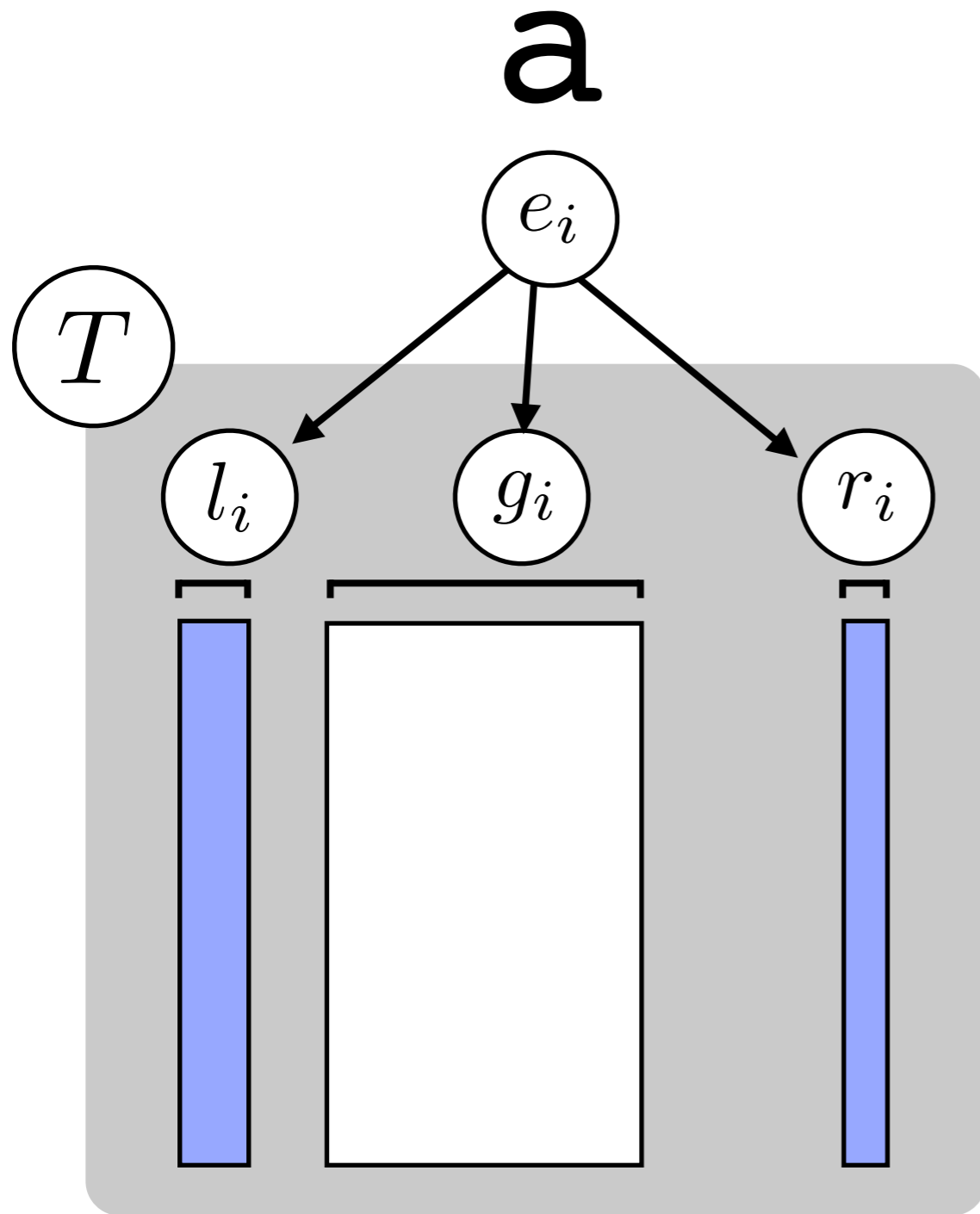
# Typesetting Model

# Typesetting Model

# Typesetting Model

a

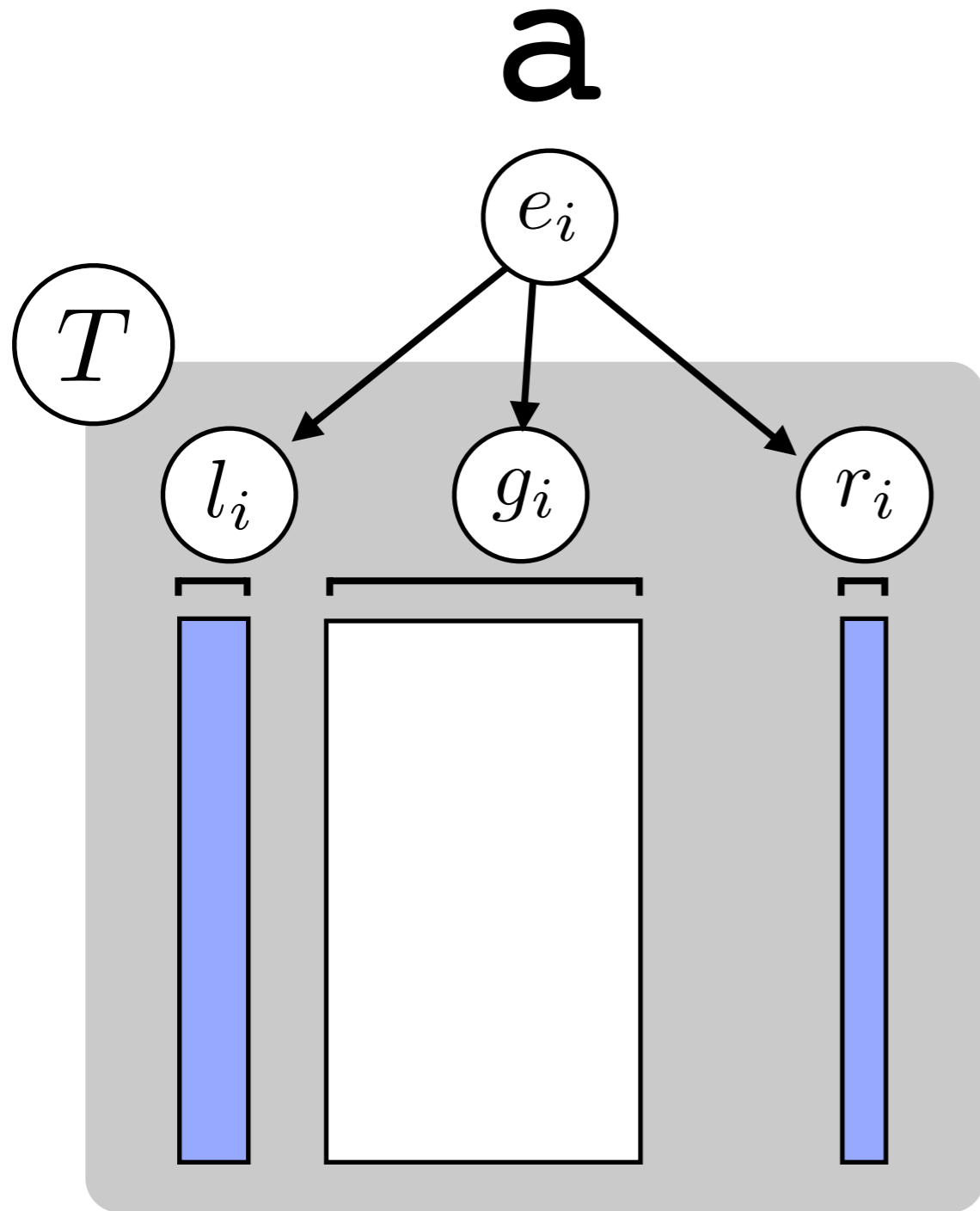$e_i$

$T$

$l_i$  $g_i$  $r_i$

$d_i$

$v_i$

Left pad width

5  1

Right pad width

1  5

Glyph box width

1  30

a

Vertical offset

Inking level

# Rendering Model

# Rendering Model

Glyph box

# Rendering Model

Glyph box
width

$g_i$

Glyph box

# Rendering Model

Glyph box width

Vertical offset

$g_i$

$v_i$

Glyph box

# Rendering Model

Glyph box width

Vertical offset

Inking level

$g_i$   $v_i$   $d_i$

Glyph box

# Rendering Model

Glyph box width

$g_i$

Vertical offset

$v_i$

Inking level

$d_i$

Glyph box

$X$

# Rendering Model

Glyph shape parameters

Glyph box width $g_i$

Vertical offset $v_i$

Inking level $d_i$

Glyph box $X$

# Rendering Model

Glyph shape parameters

Glyph box width

Vertical offset

Inking level

$g_i$

$v_i$

$d_i$

Glyph box

$X$

# Rendering Model

Glyph box width

Vertical offset

Inking level

$g_i$

$v_i$

$d_i$

Glyph shape parameters



Bernoulli pixel probs

Glyph box

$X$

# Rendering Model

Glyph shape parameters

Glyph box width

Vertical offset

Inking level

$g_i$

$v_i$

$d_i$

$X$

Glyph box

Sample pixels

Bernoulli pixel probs

# Rendering Model

Glyph box width

Vertical offset

Inking level

$g_i$

$v_i$

$d_i$

Glyph shape parameters



Glyph box

$X$

Sample pixels

Bernoulli pixel probs

# Rendering Model

Glyph shape parameters

Glyph box width

Vertical offset

Inking level

$g_i$

$v_i$

$d_i$

Glyph box

$X$

Sample pixels

Bernoulli pixel probs

# Rendering Model

Glyph shape parameters

Glyph box width

$g_i$

Vertical offset

$v_i$

Inking level

$d_i$

Bernoulli pixel probs

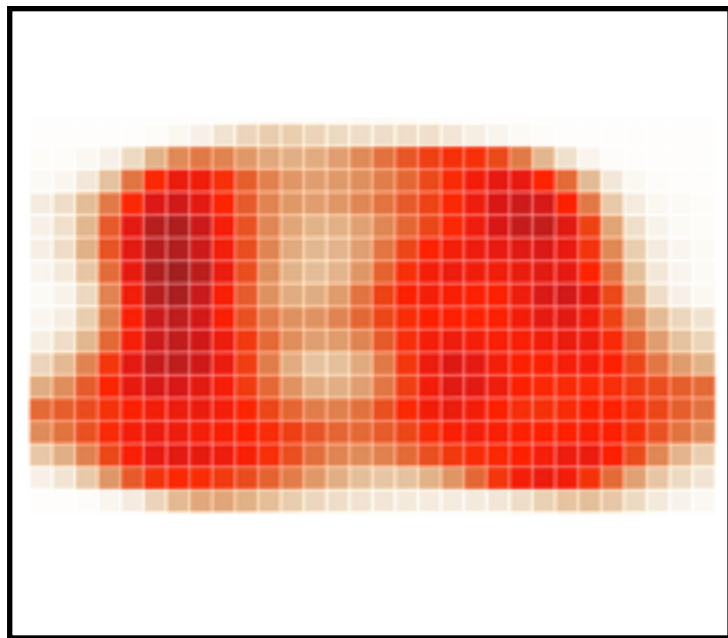Sample pixels

Glyph box

$X$

# Rendering Model

Glyph shape parameters

Glyph box width

Vertical offset

Inking level

$g_i$

$v_i$

$d_i$

$X$

Glyph box

Sample pixels

Bernoulli pixel probs

# Rendering Model

Glyph shape parameters

Glyph box width

Vertical offset

Inking level

$g_i$

$v_i$

$d_i$

Glyph box

$X$

Sample pixels

Bernoulli pixel probs

# Rendering Model

# Rendering Model

Glyph shape parameters

Glyph box width

Vertical offset

Inking level

$g_i$   $v_i$   $d_i$

Glyph box

$X$

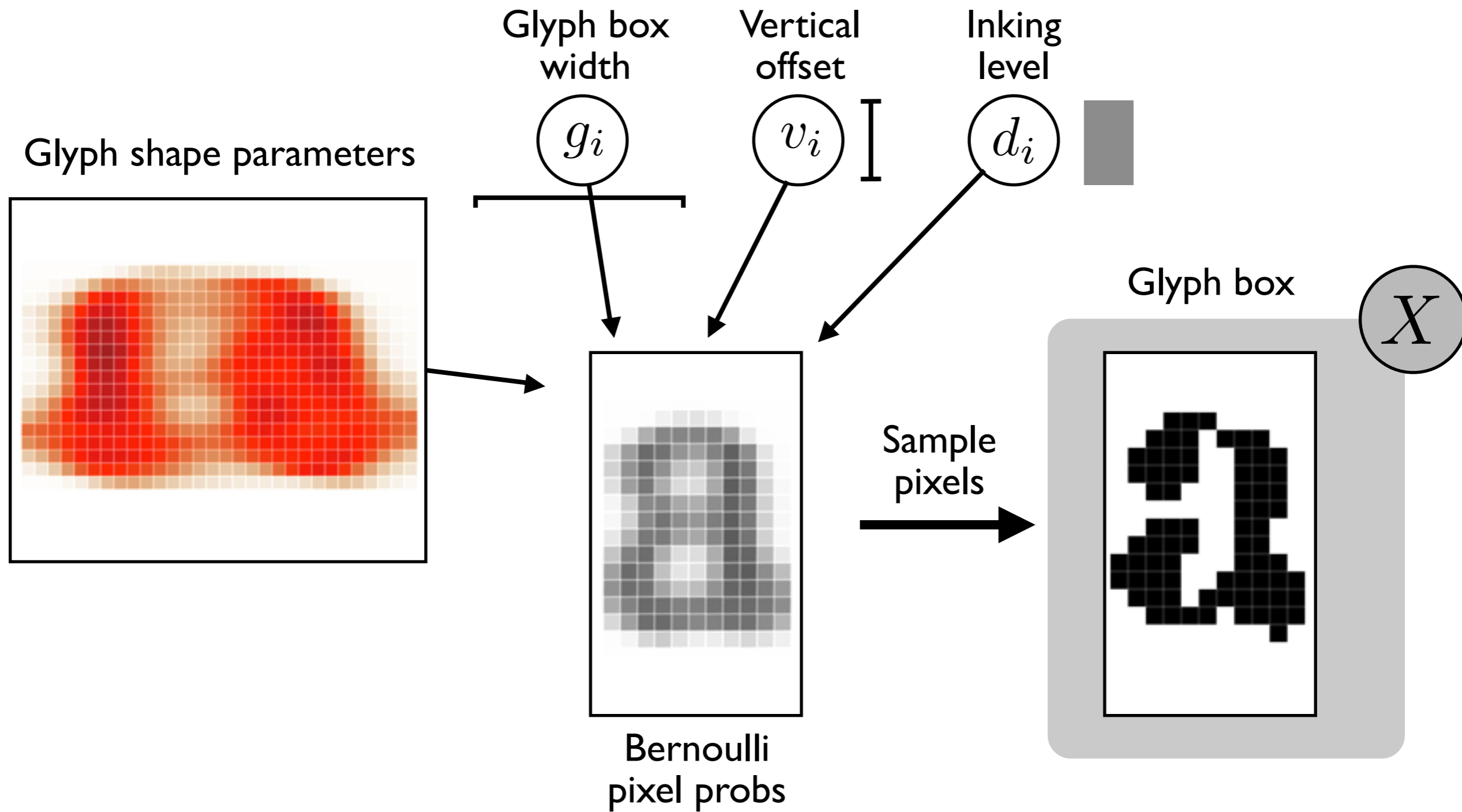Sample pixels

Bernoulli pixel probs

# Rendering Model

Glyph shape parameters

Glyph box width

Vertical offset

Inking level

$g_i$   $v_i$   $d_i$

$X$

Glyph box

Sample pixels

Bernoulli pixel probs

# Rendering Model

Glyph shape parameters

Glyph box width

$g_i$

Vertical offset

$v_i$

Inking level

$d_i$

Glyph box

$X$

Sample pixels

Bernoulli pixel probs

# Rendering Model

# Rendering Model

Glyph shape parameters

Glyph box width $g_i$

Vertical offset $v_i$

Inking level $d_i$

Glyph box $X$

Sample pixels

Bernoulli pixel probs

# Rendering Model

Glyph shape parameters

Glyph box width

Vertical offset

Inking level

$g_i$

$v_i$

$d_i$

Glyph box

$X$

Sample pixels

Bernoulli pixel probs

# Rendering Model

Glyph shape parameters

Glyph box width

Vertical offset

Inking level

$g_i$   $v_i$   $d_i$

Glyph box

$X$

Sample pixels

Bernoulli pixel probs

# Rendering Model

Glyph shape parameters

Glyph box width

Vertical offset

Inking level

$g_i$

$v_i$

$d_i$

$X$

Glyph box

Sample pixels

Bernoulli pixel probs

# Log-linear Interpolation

Glyph shape parameters $\phi$



Bernoulli pixel probs $\theta$

# Log-linear Interpolation

Glyph shape parameters $\phi$



Bernoulli pixel probs $\theta$

# Log-linear Interpolation

Glyph shape parameters $\phi$



Bernoulli pixel probs $\theta$

# Log-linear Interpolation

Glyph shape parameters $\phi$



Bernoulli pixel probs $\theta$

Glyph shape parameters $\phi$

Bernoulli pixel probs $\theta$

# Log-linear Interpolation



Interpolation weights $\alpha$

Glyph shape parameters $\phi$

Bernoulli pixel probs $\theta$

# Log-linear Interpolation



Dot product $\alpha^\top \phi$

Interpolation weights $\alpha$

Glyph shape parameters $\phi$

Bernoulli pixel probs $\theta$

# Log-linear Interpolation

Dot product $\alpha^\top \phi$

Interpolation weights $\alpha$

Glyph shape parameters $\phi$

Apply logistic

Bernoulli pixel probs $\theta$

# Log-linear Interpolation

Interpolation weights $\alpha$

Glyph shape parameters $\phi$

Bernoulli pixel probs $\theta$

# Log-linear Interpolation



Interpolation weights $\alpha$

Glyph shape parameters $\phi$

Bernoulli pixel probs $\theta$

# Log-linear Interpolation



Interpolation weights $\alpha$

Glyph shape parameters $\phi$

Bernoulli pixel probs $\theta$

# Log-linear Interpolation

Interpolation weights $\alpha_j$

Glyph shape parameters $\phi$

Bernoulli pixel probs $\theta$

$j$

# Log-linear Interpolation

$$\theta_j \propto$$

Interpolation weights $\alpha_j$

Glyph shape parameters $\phi$

Bernoulli pixel probs $\theta$

$j$

# Log-linear Interpolation

$$\theta_j \propto \exp[\alpha_j^\top \phi]$$

Interpolation weights $\alpha_j$

Glyph shape parameters $\phi$

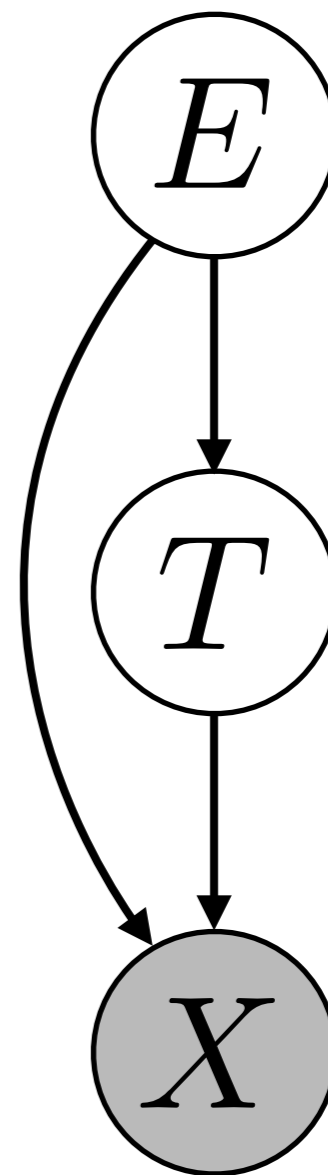Bernoulli pixel probs $\theta$

$j$

- Learn font parameters using EM

- Learn font parameters using EM

- Learn font parameters using EM

- Learn font parameters using EM
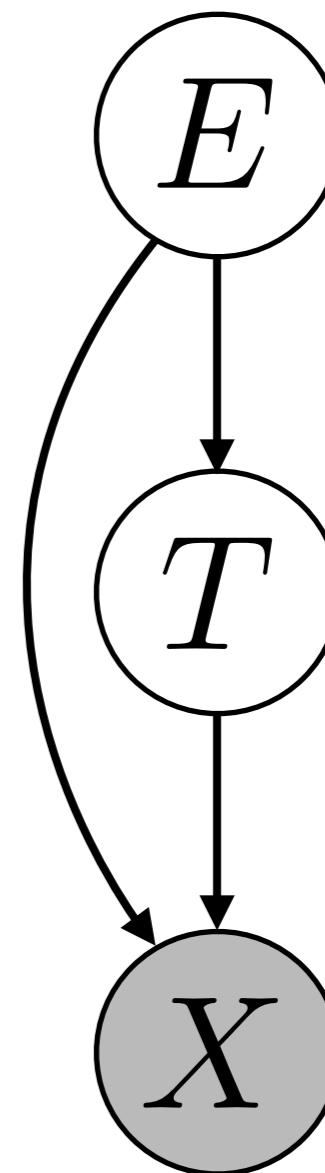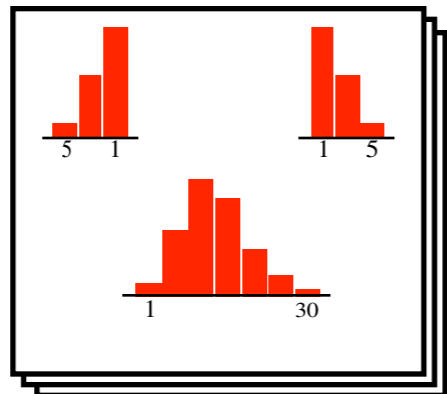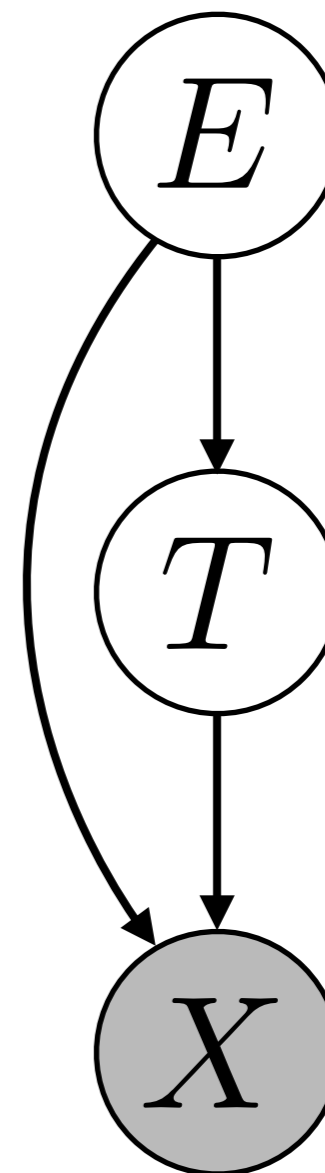


- Initialize font parameters with mixtures of modern fonts

# Learning and Inference

- Learn font parameters using EM

- Initialize font parameters with mixtures of modern fonts

- Semi-Markov DP to compute expectations
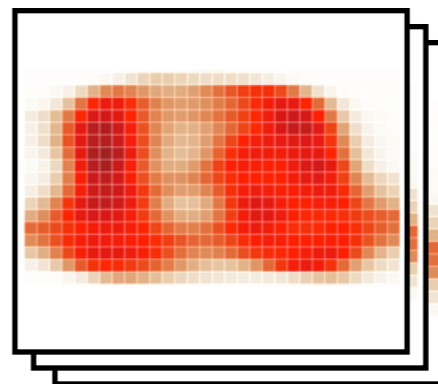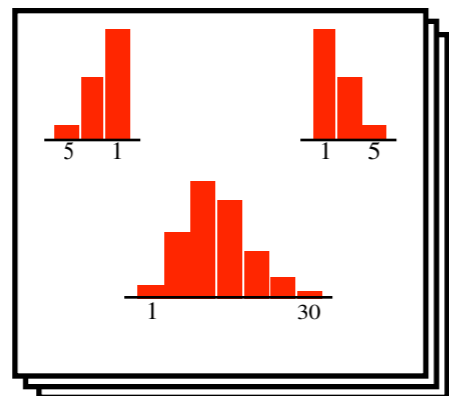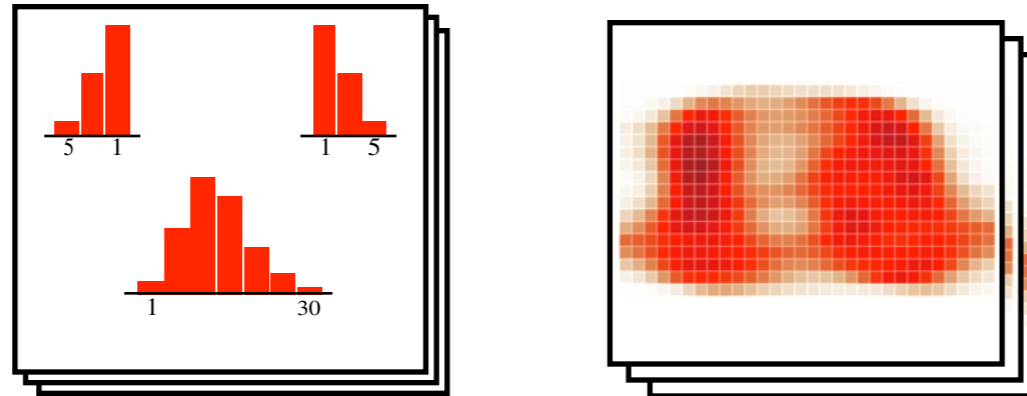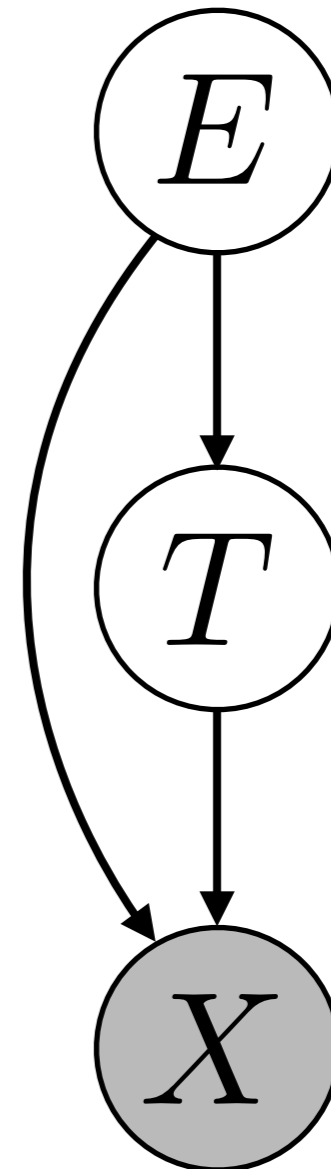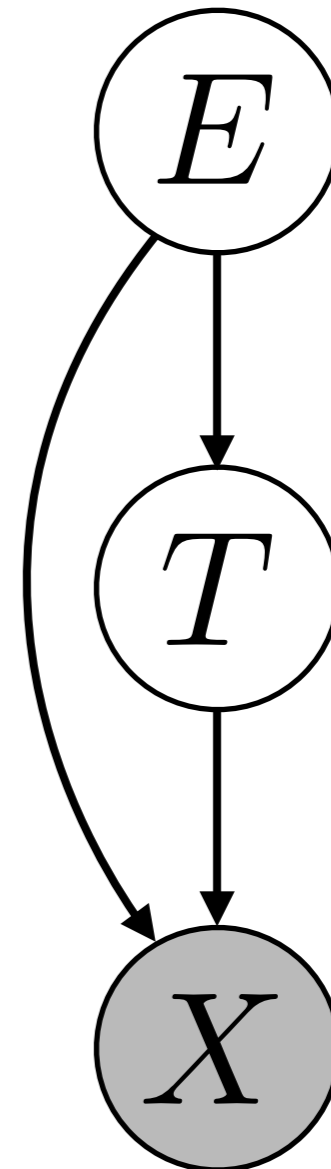
# Learning and Inference

- Learn font parameters using EM



- Initialize font parameters with mixtures of modern fonts

- Semi-Markov DP to compute expectations

- Efficient inference using a coarse-to-fine approach

how the murderers came to

how the murderers came to

# System Output Example



how the murderers came to

# System Output Example



```
how the murderers came to
```

# System Output Example

how the murderers came to

# System Output Example

how the murderers came to

taken ill and taken away—I remember

taken ill and taken away—I remember

# System Output Example



`taken ill and taken  away -- I remember`
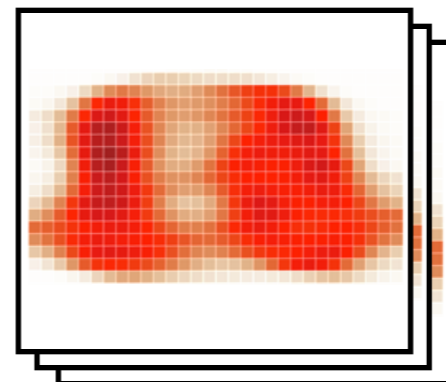
# System Output Example



taken ill and taken  away -- I remember

# System Output Example



taken ill and taken  away -- I remember

# System Output Example



taken ill and taken  away -- I remember

Test data

## Test data

- Old Bailey (1715-1905)

## Test data

- Old Bailey (1715-1905)

  20 images, 30 lines each

## Test data

- Old Bailey (1715-1905)

    20 images, 30 lines each

- Trove (1803-1893)

## Test data

- Old Bailey (1715-1905)

  20 images, 30 lines each

- Trove (1803-1893)

  10 images, 30 lines each

Experiments

## Test data

- Old Bailey (1715-1905)

  20 images, 30 lines each

- Trove (1803-1893)

  10 images, 30 lines each

## Baselines

# Experiments

## Test data

- Old Bailey (1715-1905)

    20 images, 30 lines each

- Trove (1803-1893)

    10 images, 30 lines each

## Baselines

- Google Tesseract

# Experiments

## Test data

- Old Bailey (1715-1905)

  20 images, 30 lines each

- Trove (1803-1893)

  10 images, 30 lines each

## Baselines

- Google Tesseract

- ABBYY FineReader 11

# Experiments

## Test data

- Old Bailey (1715-1905)

    20 images, 30 lines each

- Trove (1803-1893)

    10 images, 30 lines each

## Baselines

- Google Tesseract

- ABBYY FineReader 11

## Language models

# Experiments

## Test data

- Old Bailey (1715-1905)

    20 images, 30 lines each

- Trove (1803-1893)

    10 images, 30 lines each

## Baselines

- Google Tesseract

- ABBYY FineReader 11

## Language models

- New York Times

# Experiments

## Test data

- Old Bailey (1715-1905)

  20 images, 30 lines each

- Trove (1803-1893)

  10 images, 30 lines each

## Language models

- New York Times

  34M words NYT Gigaword

## Baselines

- Google Tesseract

- ABBYY FineReader 11

# Experiments

## Test data

- Old Bailey (1715-1905)

  20 images, 30 lines each

- Trove (1803-1893)

  10 images, 30 lines each

## Language models

- New York Times

  34M words NYT Gigaword

- Old Bailey

## Baselines

- Google Tesseract

- ABBYY FineReader 11

# Experiments

## Test data

- Old Bailey (1715-1905)

  20 images, 30 lines each

- Trove (1803-1893)

  10 images, 30 lines each

## Language models

- New York Times

  34M words NYT Gigaword

- Old Bailey

  32M words manually transcribed

## Baselines

- Google Tesseract

- ABBYY FineReader 11

Old Bailey Court Proceedings
(1715-1905)

# Results

## Old Bailey Court Proceedings
## (1715-1905)

# Results

## Old Bailey Court Proceedings (1715-1905)

# Results



Old Bailey Court Proceedings
(1715-1905)

[Berg-Kirkpatrick et al. 2013]

# Results



Old Bailey Court Proceedings
(1715-1905)

[Berg-Kirkpatrick et al. 2013]

# Results

## Trove Historical Newspapers
## (1803-1893)

# Results

## Trove Historical Newspapers
## (1803-1893)

# Results

## Trove Historical Newspapers (1803-1893)

# Results

## Trove Historical Newspapers
## (1803-1893)

# Results

## Trove Historical Newspapers
## (1803-1893)



[Berg-Kirkpatrick et al. 2014]

## Google Tesseract

and Ch': priftmer anhc bar. Jacob Lazarus and his
IHP1 uh: prifoner. were both together when!
rcccivcd lhczn. I fold eievén pair of than
for xiirce guincas, and dclivcrcd the rcll'l:.in-
d:r hack lo :11: prifuner. 1 fold ftvcn pairof
filk to Mark Simpcr : nncpuir of mixcd. and.
mo pair of Ifircad to lhz: foolnun, and on:
pair of zhrzad to lh: barber. '
  Q: What is the foolmarfs name?
  Fraum Mgfzr.   I dun': know.
  Hairy Hzrvir. l was flandingar the Camp
Icr waizin far the thcrrilfs ufliceruo employ
in: : Mo 3': daughter came for me to 0 am!
take the prifoncr. 1 Wm! to |hc Old aailcy

# Transcription

## Google Tesseract

and Ch': priftmer anhc bar. Jacob Lazarus and his
IHP1 uh: prifoner. were both together when!
rcccivcd lhczn. I fold eievén pair of than
for xiirce guincas, and dclivcrcd the rcll'l:.in-
d:r hack lo :11: prifuner. 1 fold ftvcn pairof
filk to Mark Simpcr :   nncpuir of mixcd. and.
mo pair of Ifircad to lhz: foolnun, and on:
pair of zhrzad to lh: barber. '
   Q: What is the foolmarfs name?
   Fraum Mgfzr.   I dun': know.
   Hairy Hzrvir. l was flandingar the Camp
Icr waizin far the thcrrilfs ufliceruo employ
in: : Mo 3': daughter came for me to 0 am!
take the prifoncr. 1 Wm! to |hc Old aailcy

## Ocular

the prisoner at the bar. Jacob Lazarus and his
wife, the prisoners were both together when I
received them. I sold eleven pair of them
for three guineas, and delivered the remain-
der back to the prisoner. I sold, seven pair of
silk to Mark Simpert one pair of mixed, and
two pair of thread to the footman, and one
pair of thread to the barber,
   Ms. What in the footman's name?
   Franco Asyut,  I don't know-
   Nearly Norris. I was standing at the Comp-
ter waiting for the sherrill's officers to employ
me a Moses's daughter came for me to go and
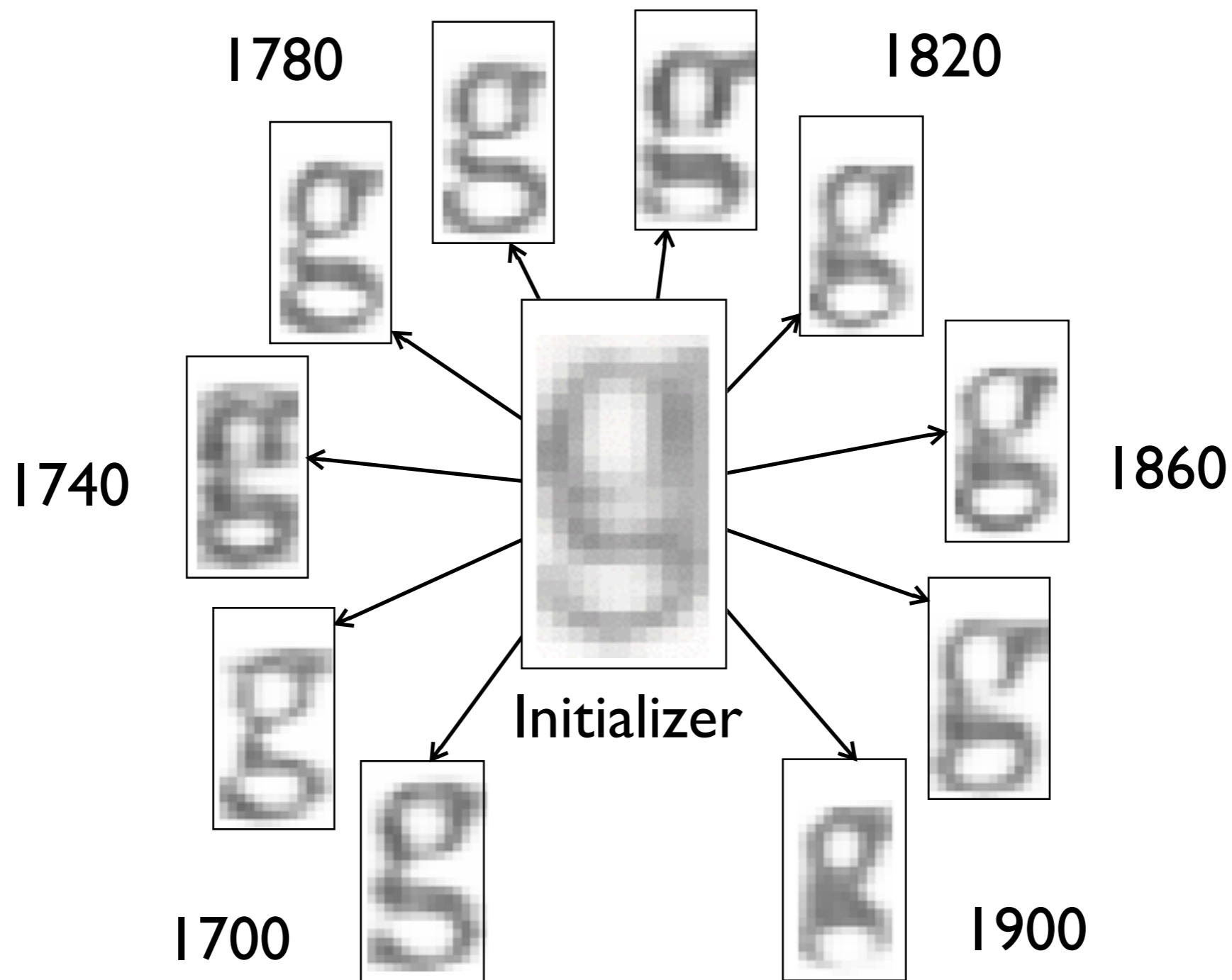take the prisoner. I went to the Old Bailey

Initializer

Initializer

Initializer

# Learned Fonts



1780

1820

1740

1860

1700

1900

Initializer

# Unobserved Pixels

# Unobserved Pixels

# Unobserved Pixels

# Unobserved Pixels

- Unsupervised font learning yields state-of-the-art results on documents where font is unknown

# Conclusion

- Unsupervised font learning yields state-of-the-art results on documents where font is unknown

- Generatively modeling sources of noise specific to printing-press era documents is effective

# Conclusion

- Unsupervised font learning yields state-of-the-art results on documents where font is unknown

- Generatively modeling sources of noise specific to printing-press era documents is effective

- Ocular available as a downloadable tool:
  nlp.cs.berkeley.edu/ocular.shtml

# Thanks!