

# Structure Compilation: Trading Structure for Features

ICML 2008

Helsinki, Finland

July 8, 2008

Percy Liang  
UC Berkeley

Hal Daumé  
U of Utah

Dan Klein  
UC Berkeley

## Structured prediction:

$y$ : DT — NNP — NNP — VBD  
 $x$ : The European Commision agreed

Part-of-speech tagging (POS)

## Structured prediction:

$y$ : DT — NNP — NNP — VBD  
 $x$ : The European Commission agreed  
Part-of-speech tagging (POS)

$y$ : O — B-ORG — I-ORG — O  
 $x$ : The European Commission agreed  
Named-entity recognition (NER)

## Structured prediction:

$y$ : DT — NNP — NNP — VBD  
 $x$ : The European Commission agreed

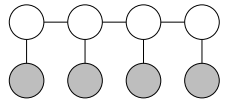
Part-of-speech tagging (POS)

$y$ : O — B-ORG — I-ORG — O  
 $x$ : The European Commission agreed

Named-entity recognition (NER)

## Current methods:

Structured models: **accurate** but **slow**



conditional random fields (CRFs) with loopy graphs, large tag sets

## Structured prediction:

$y$ : DT — NNP — NNP — VBD  
 $x$ : The European Commission agreed

Part-of-speech tagging (POS)

$y$ : O — B-ORG — I-ORG — O  
 $x$ : The European Commission agreed

Named-entity recognition (NER)

## Current methods:

Structured models: **accurate** but **slow**

 conditional random fields (CRFs) with loopy graphs, large tag sets

Independent models: **less accurate** but **fast**

 independent logistic regressions (ILRs)

## Structured prediction:

$y$ : DT — NNP — NNP — VBD  
 $x$ : The European Commission agreed

Part-of-speech tagging (POS)

$y$ : O — B-ORG — I-ORG — O  
 $x$ : The European Commission agreed

Named-entity recognition (NER)

## Current methods:

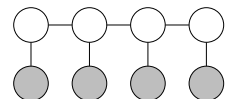
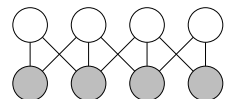
Structured models: accurate but slow

 conditional random fields (CRFs) with loopy graphs, large tag sets

Independent models: less accurate but fast

 independent logistic regressions (ILRs)

## Our goal:

  $\xrightarrow[\text{predictive power}]{\text{transfer}}$   **accurate and fast at test time**

## Structured prediction:

$y$ : DT — NNP — NNP — VBD  
 $x$ : The European Commission agreed

Part-of-speech tagging (POS)

$y$ : O — B-ORG — I-ORG — O  
 $x$ : The European Commission agreed

Named-entity recognition (NER)

## Current methods:

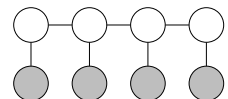
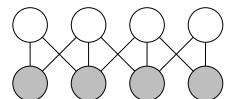
Structured models: accurate but slow

 conditional random fields (CRFs) with loopy graphs, large tag sets

Independent models: less accurate but fast

 independent logistic regressions (ILRs)

## Our goal:

  $\xrightarrow[\text{predictive power}]{\text{transfer}}$   **accurate and fast at test time**

Questions: are independent models...

- ...expressive enough (**approximation error**)?

## Structured prediction:

$y$ : DT — NNP — NNP — VBD  
 $x$ : The European Commission agreed

Part-of-speech tagging (POS)

$y$ : O — B-ORG — I-ORG — O  
 $x$ : The European Commission agreed

Named-entity recognition (NER)

## Current methods:

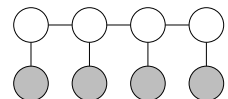
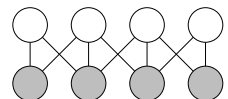
Structured models: accurate but slow

 conditional random fields (CRFs) with loopy graphs, large tag sets

Independent models: less accurate but fast

 independent logistic regressions (ILRs)

## Our goal:

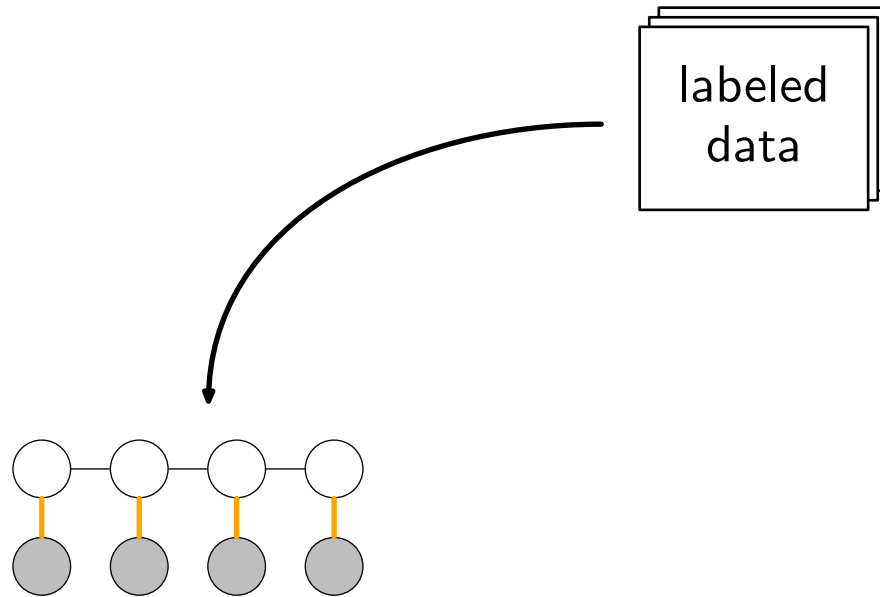
  $\xrightarrow[\text{predictive power}]{\text{transfer}}$   **accurate and fast at test time**

## Questions: are independent models...

- ...expressive enough (**approximation error**)?
- ...easy to learn (**estimation error**)?



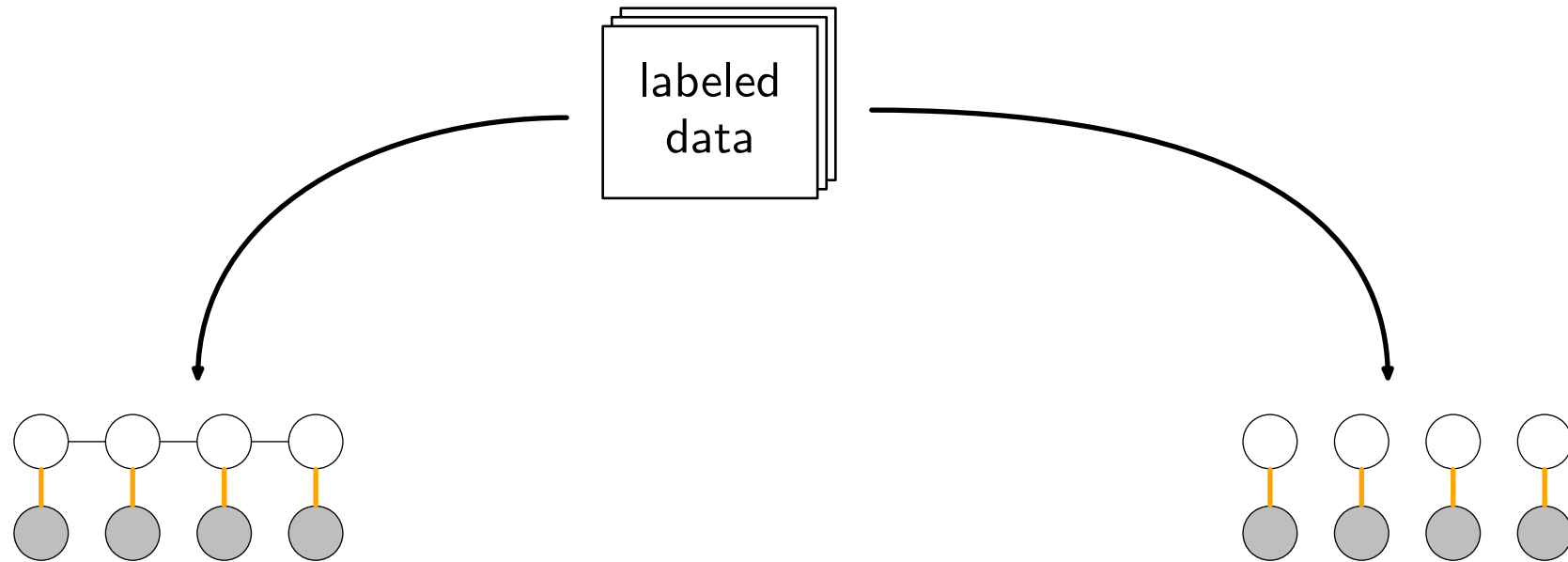
# Some empirical motivation



CRF( $f_1$ ) POS: 95.0%  
NER: 75.3%

$f_1$ : words/prefixes/suffixes/forms

# Some empirical motivation

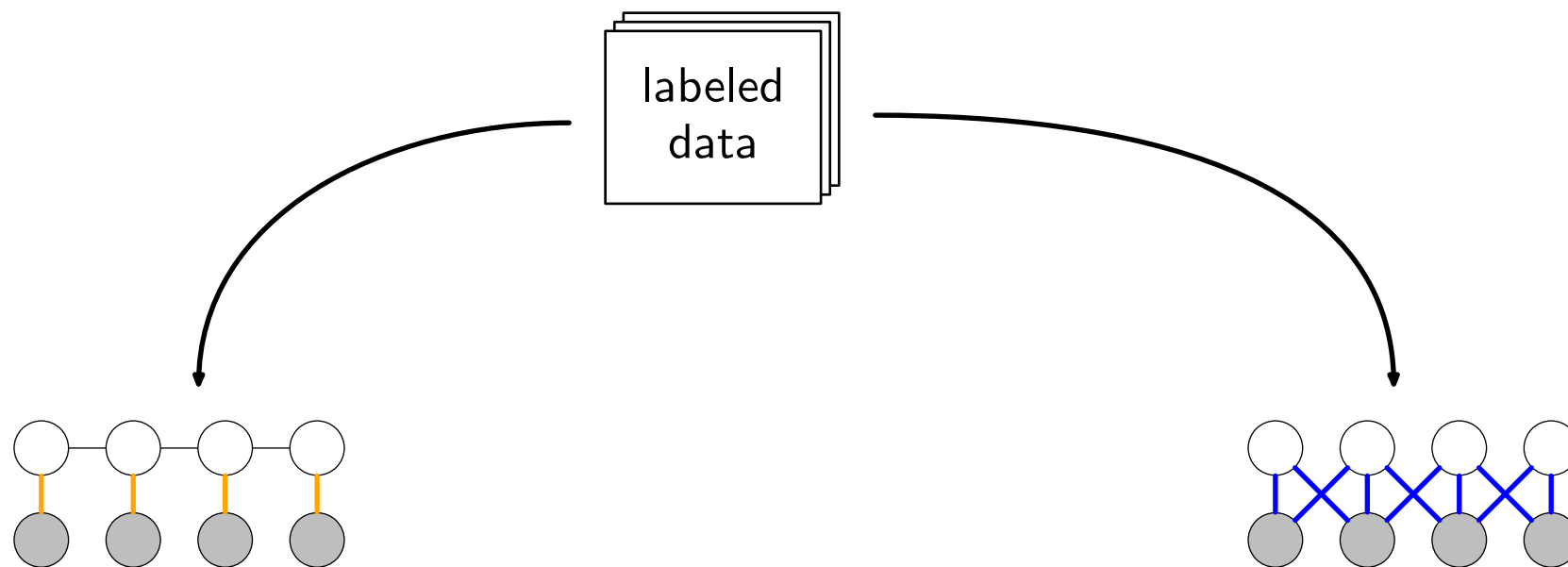


CRF( $f_1$ ) POS: 95.0%  
NER: 75.3%

ILR( $f_1$ ) POS: 91.7%  
NER: 69.1%

$f_1$ : words/prefixes/suffixes/forms

# Some empirical motivation



CRF( $f_1$ ) POS: 95.0%  
NER: 75.3%

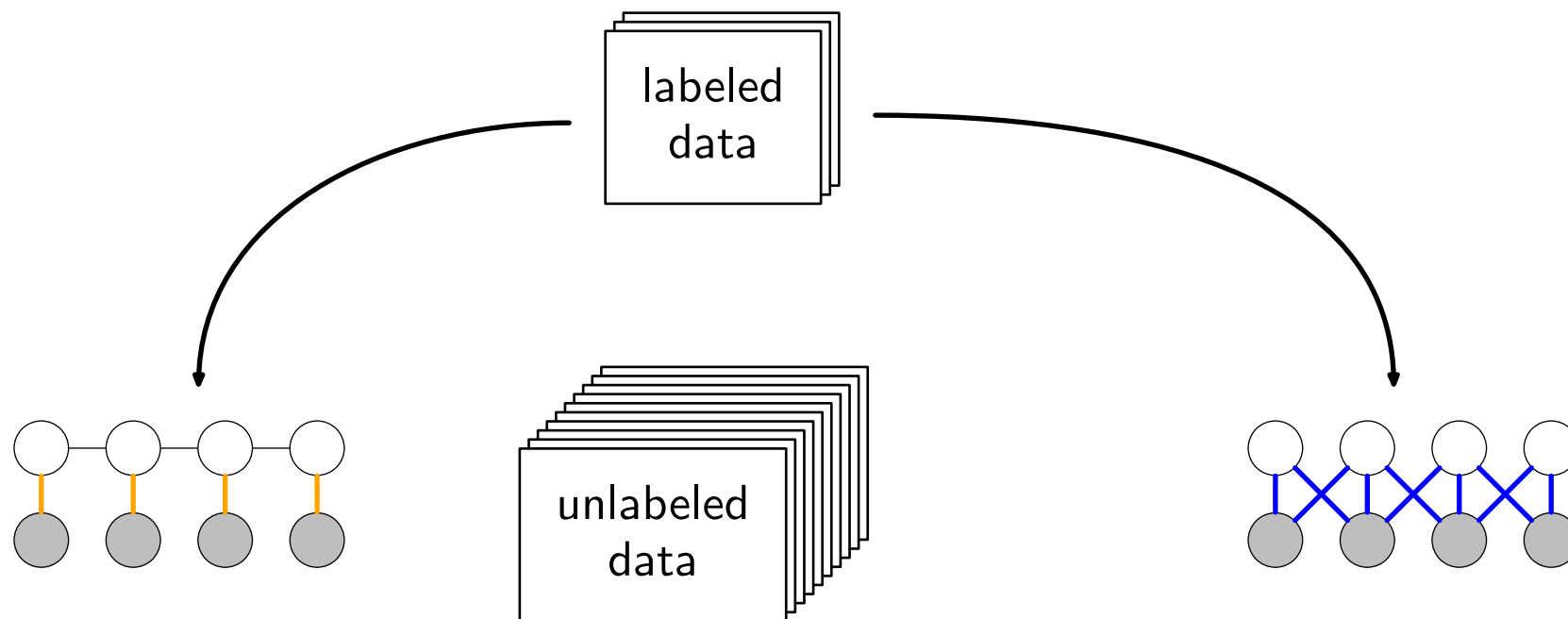
$f_1$ : words/prefixes/suffixes/forms

$f_2$ :  $f_1$  applied to larger radius

ILR( $f_1$ ) POS: 91.7%  
NER: 69.1%

ILR( $f_2$ ) POS: 94.4%  
NER: 66.2%

# Some empirical motivation



CRF( $f_1$ ) POS: 95.0%  
NER: 75.3%

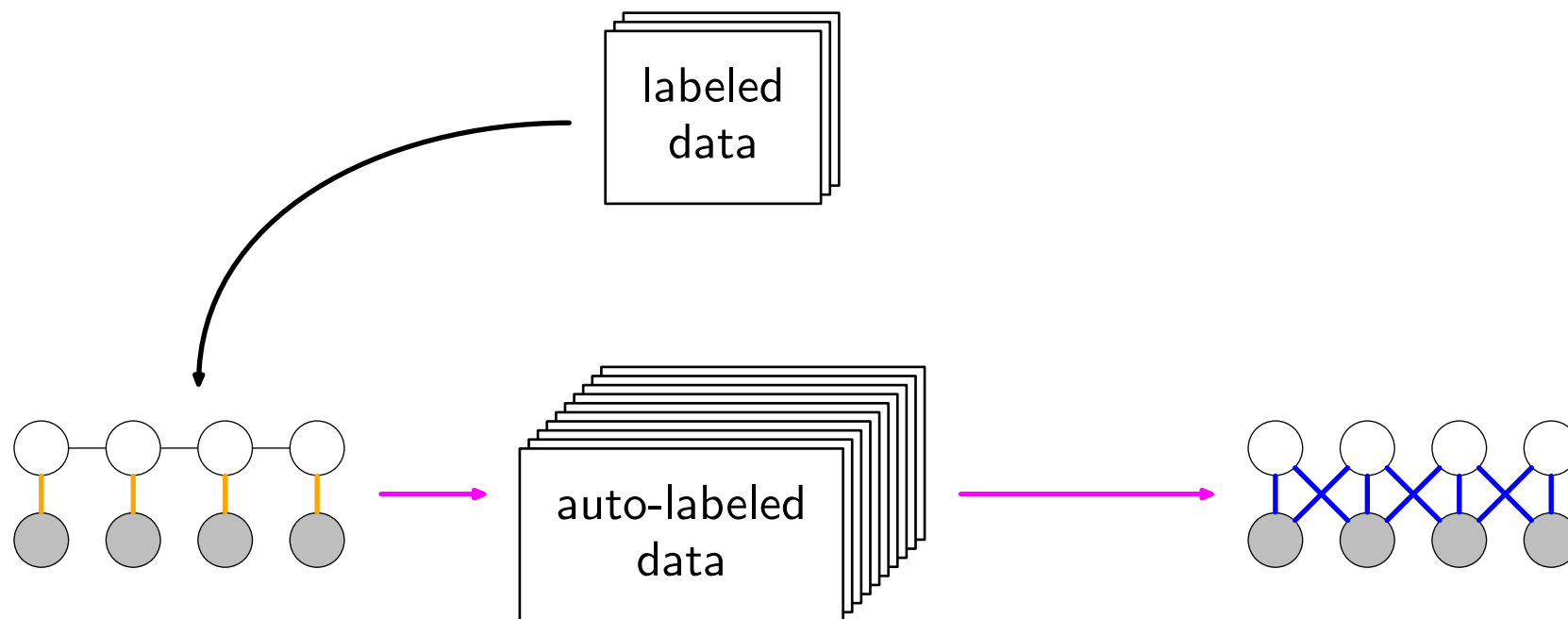
$f_1$ : words/prefixes/suffixes/forms

$f_2$ :  $f_1$  applied to larger radius

ILR( $f_1$ ) POS: 91.7%  
NER: 69.1%

ILR( $f_2$ ) POS: 94.4%  
NER: 66.2%

# Some empirical motivation



CRF( $f_1$ ) POS: 95.0%  
NER: 75.3%

$f_1$ : words/prefixes/suffixes/forms

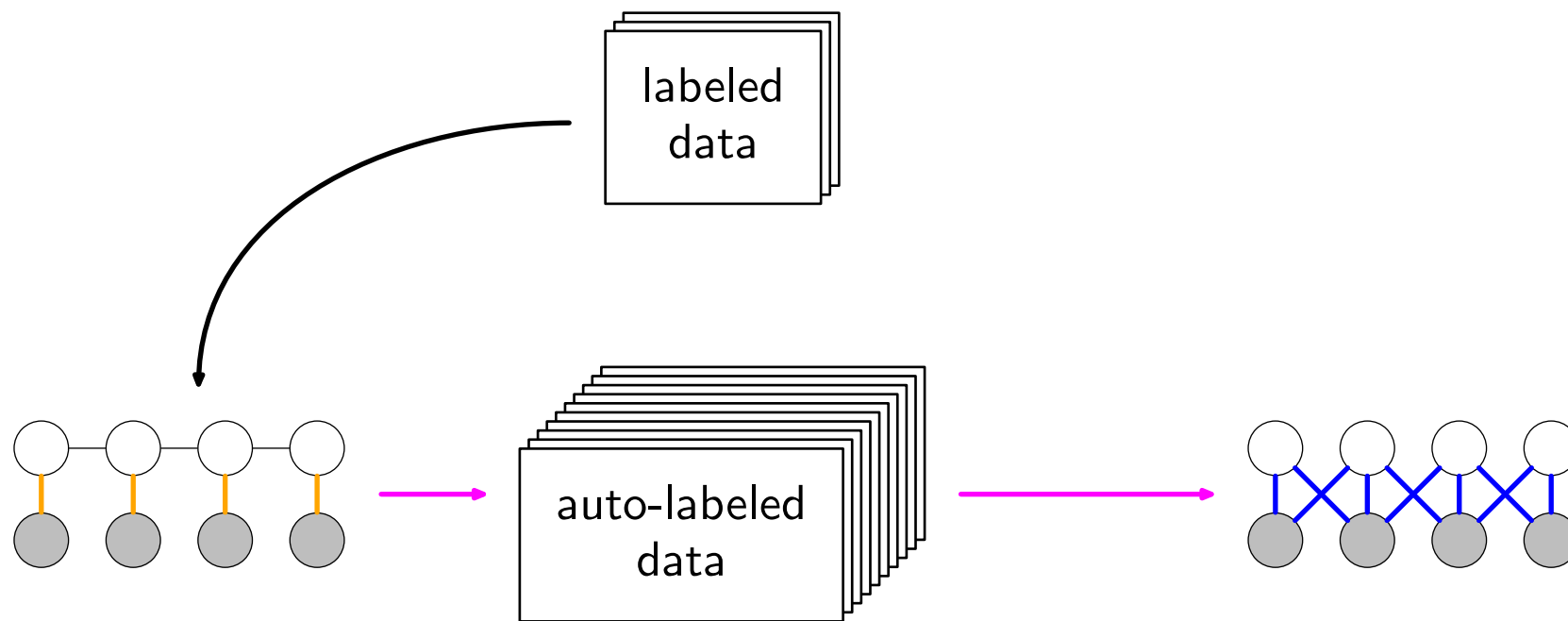
$f_2$ :  $f_1$  applied to larger radius

ILR( $f_1$ ) POS: 91.7%  
NER: 69.1%

ILR( $f_2$ ) POS: 94.4%  
NER: 66.2%

ILR( $f_2$ ) [compiled] POS: 95.0%  
NER: 72.7%

# Some empirical motivation



CRF( $f_1$ ) POS: 95.0%  
NER: 75.3%

$f_1$ : words/prefixes/suffixes/forms

$f_2$ :  $f_1$  applied to larger radius

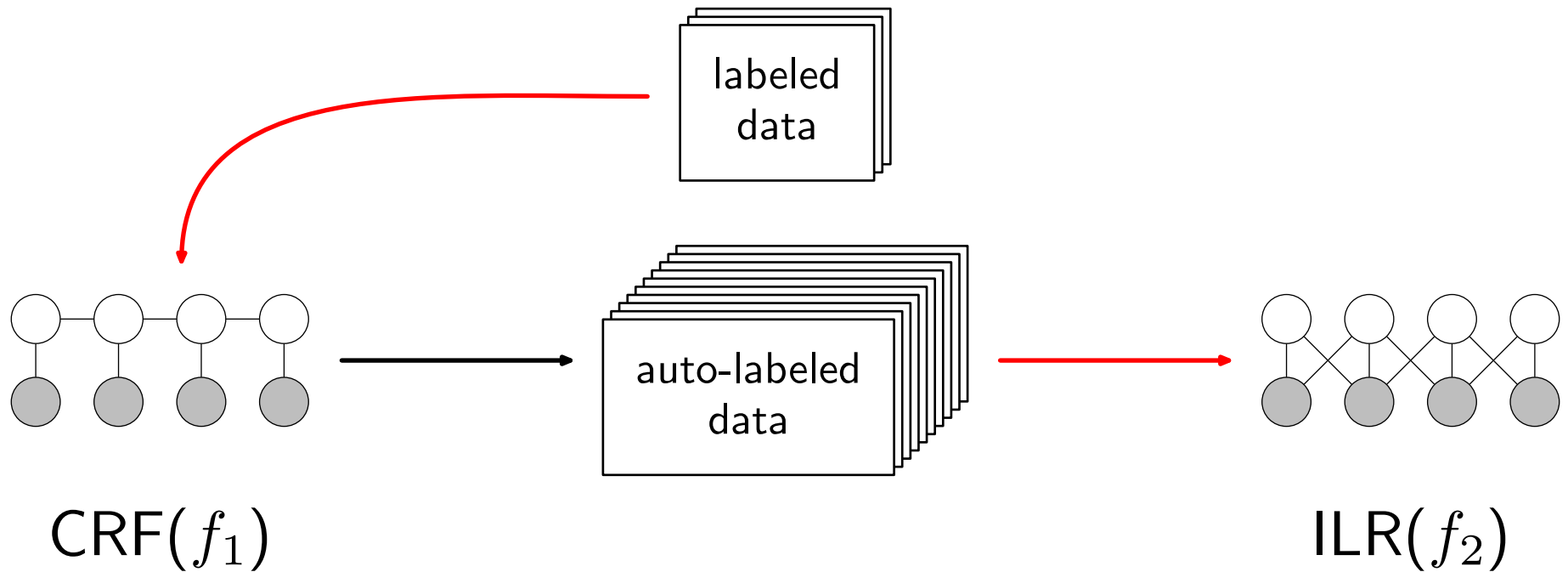
ILR( $f_1$ ) POS: 91.7%  
NER: 69.1%

ILR( $f_2$ ) POS: 94.4%  
NER: 66.2%

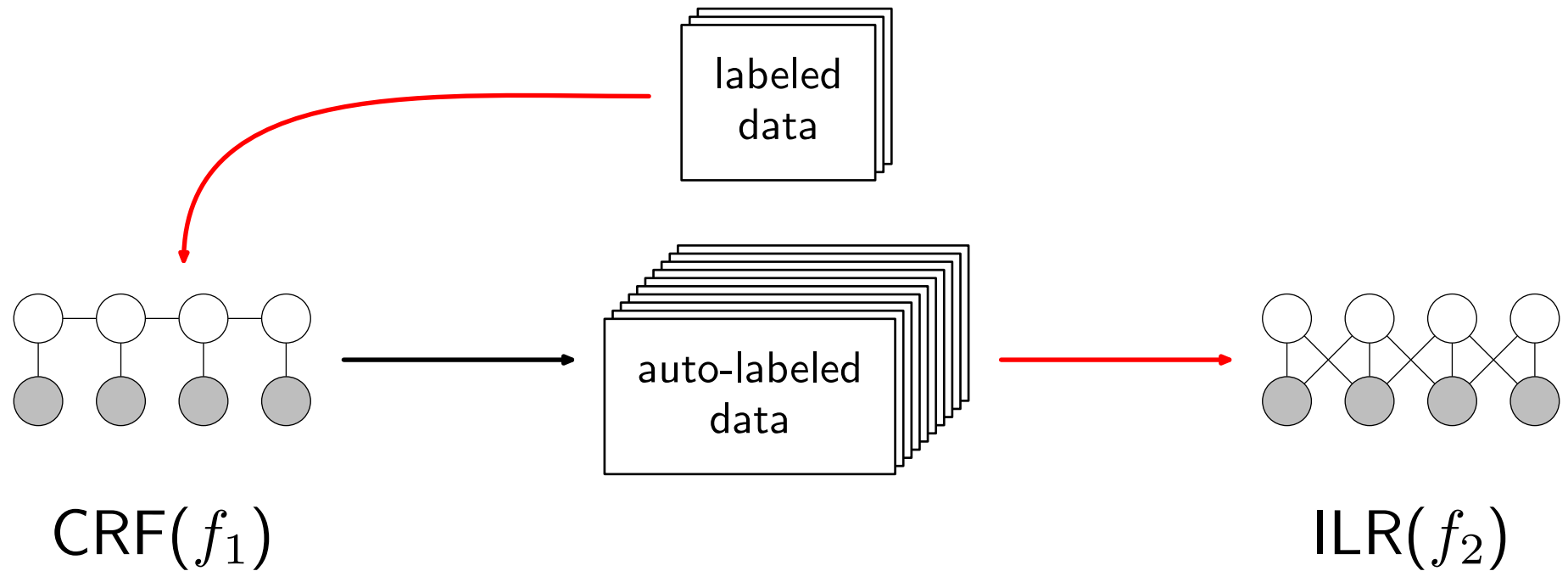
ILR( $f_2$ ) [compiled] POS: 95.0%  
NER: 72.7%

Structure compilation: reduces the gap between the ILR and CRF

# Analysis of structure compilation



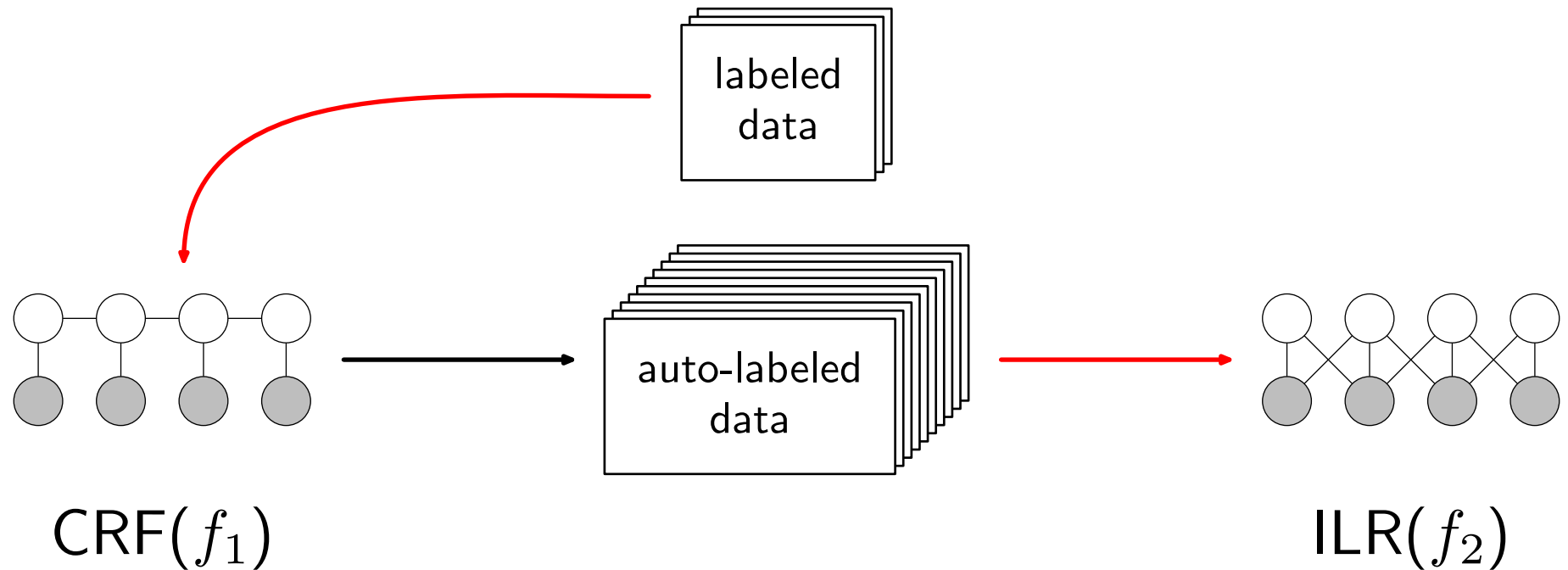
# Analysis of structure compilation



**Goal:** analyze risk of final compiled  $ILR(f_2)$



# Analysis of structure compilation

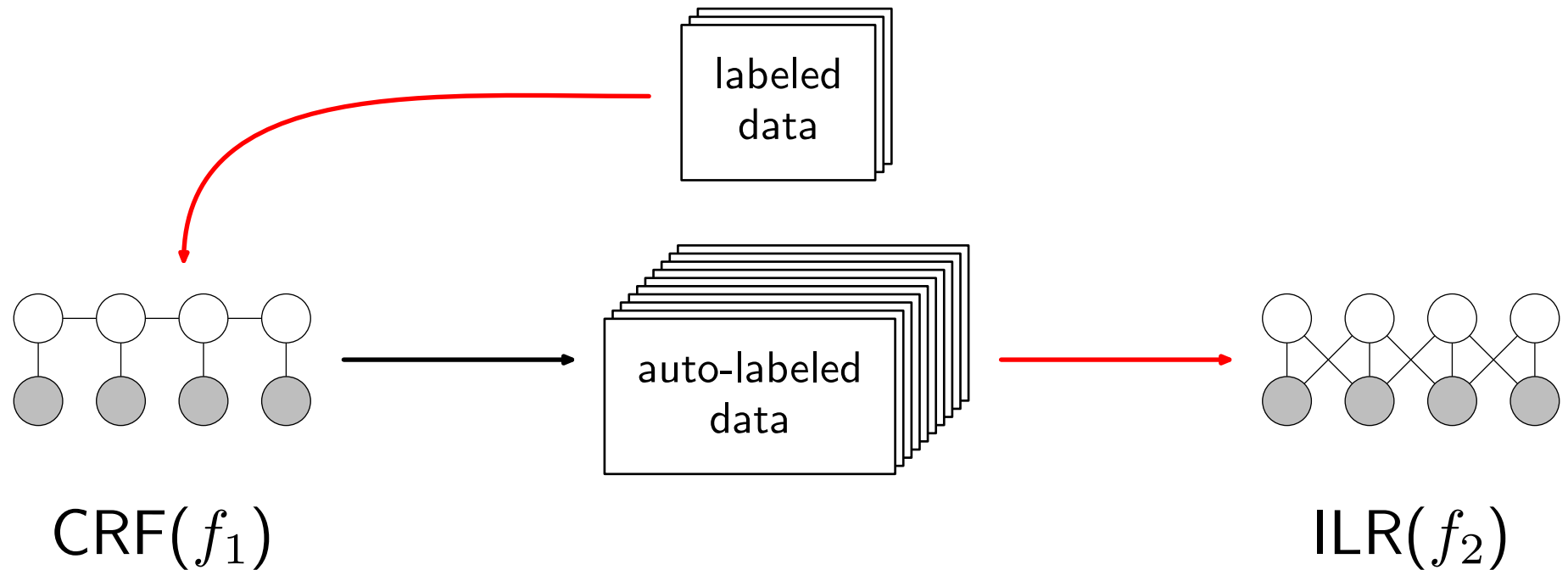


Goal: analyze risk of final compiled  $ILR(f_2)$

Decomposition of errors:

**Approximation error:** best loss of model (with infinite data)

# Analysis of structure compilation



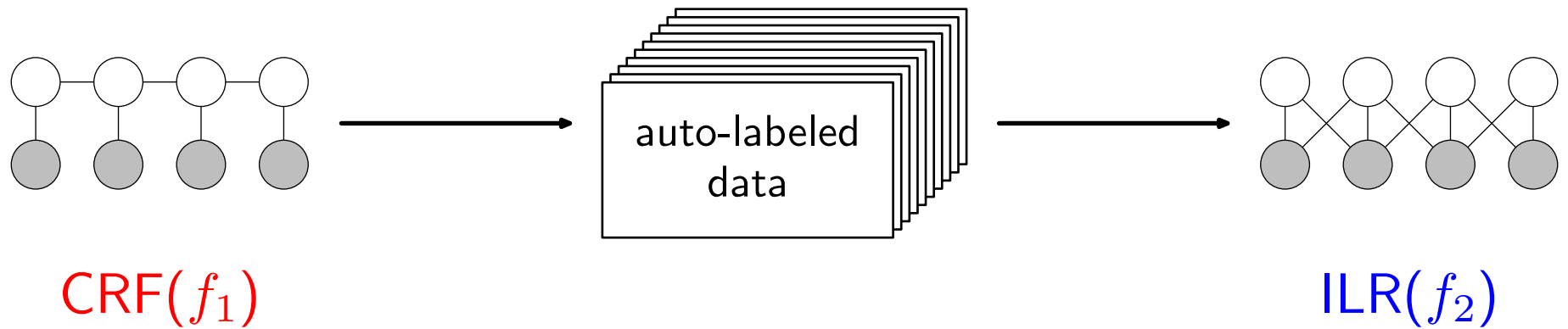
Goal: analyze risk of final compiled  $ILR(f_2)$

Decomposition of errors:

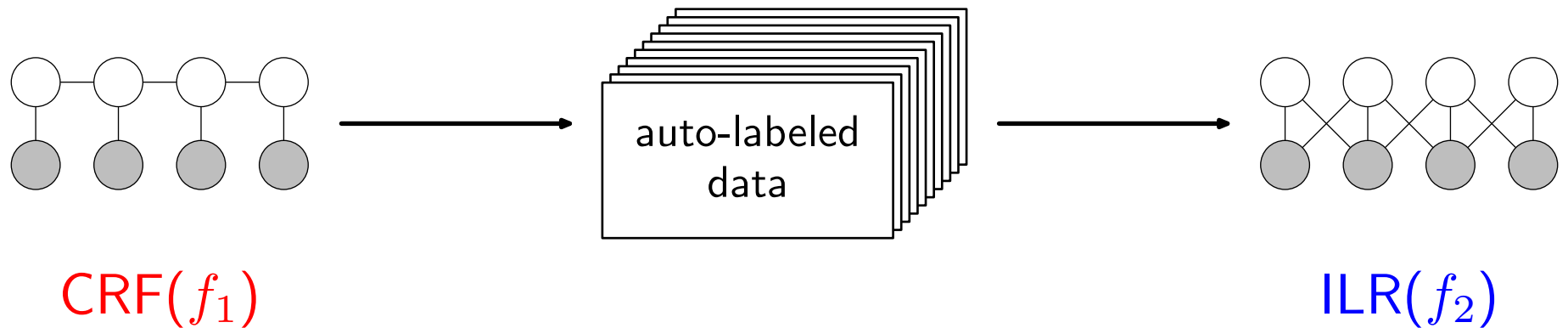
**Approximation error:** best loss of model (with infinite data)

**Estimation error:** suboptimality due to finite data

# Approximation/estimation errors for ILR

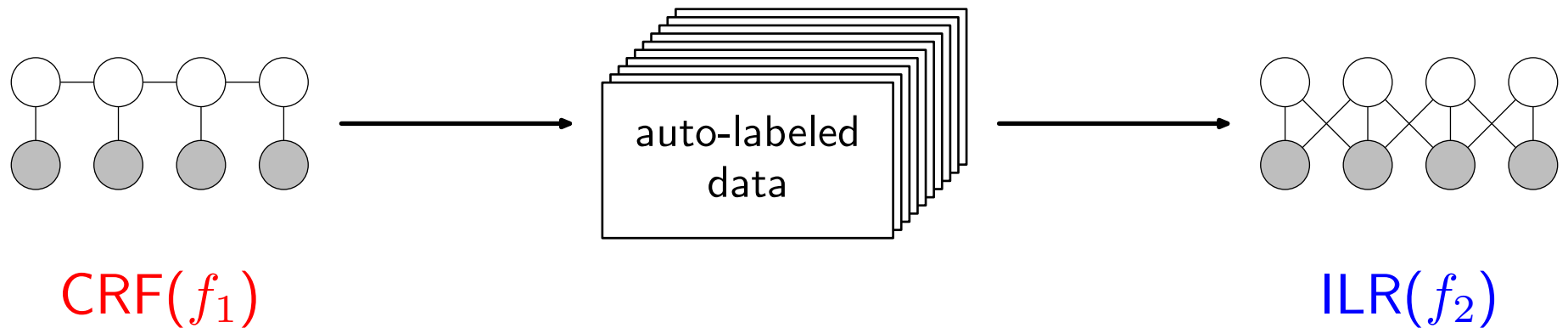


# Approximation/estimation errors for ILR



$p_C$ :  $CRF(f_1)$  trained on labeled data

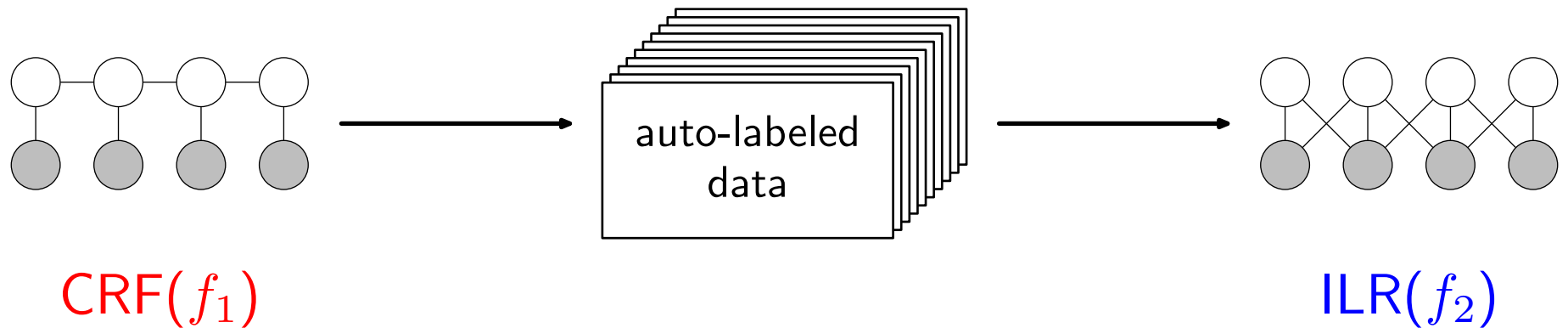
# Approximation/estimation errors for ILR



$p_C$ :  $CRF(f_1)$  trained on labeled data

$p_I$ :  $ILR(f_2)$  trained on  $m$  auto-labeled examples

# Approximation/estimation errors for ILR

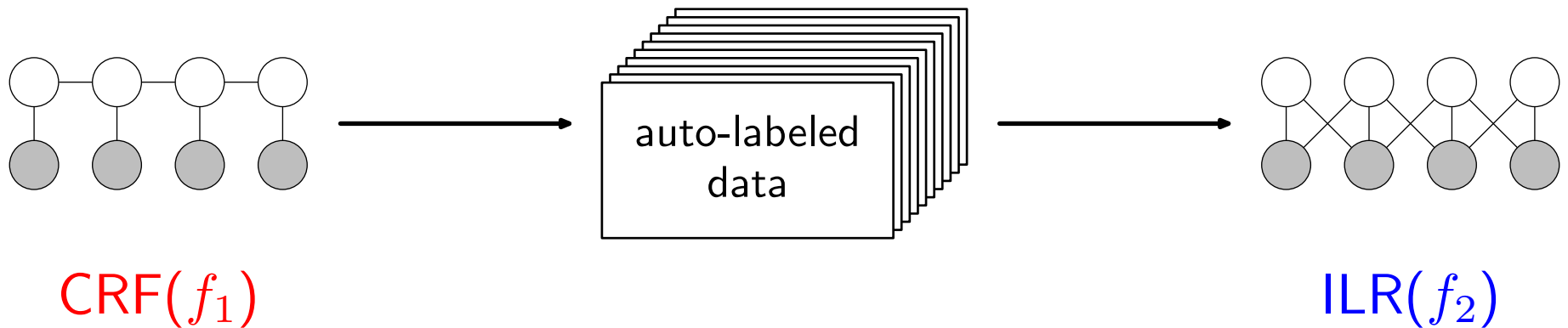


$p_C$ :  $CRF(f_1)$  trained on labeled data

$p_I$ :  $ILR(f_2)$  trained on  $m$  auto-labeled examples

$p_{I^*}$ :  $ILR(f_2)$  trained on infinite auto-labeled data

# Approximation/estimation errors for ILR



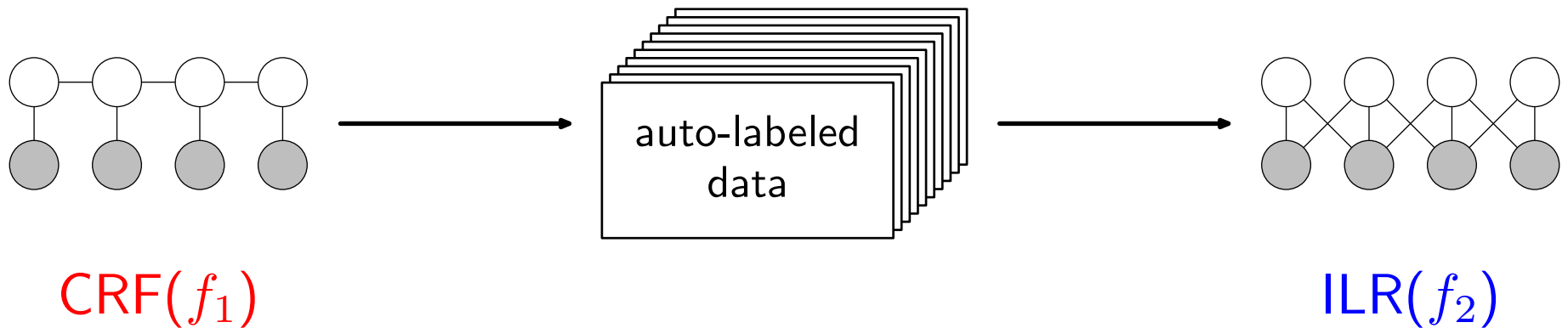
$p_C$ :  $CRF(f_1)$  trained on labeled data

$p_I$ :  $ILR(f_2)$  trained on  $m$  auto-labeled examples

$p_{I^*}$ :  $ILR(f_2)$  trained on infinite auto-labeled data

$$\text{KL} (p_C || p_I) = \underbrace{\text{KL} (p_C || p_{I^*})}_{\text{approx. error}} + \underbrace{(\text{KL} (p_C || p_I) - \text{KL} (p_C || p_{I^*}))}_{\text{estimation error}}$$

# Approximation/estimation errors for ILR



$p_C$ :  $CRF(f_1)$  trained on labeled data

$p_I$ :  $ILR(f_2)$  trained on  $m$  auto-labeled examples

$p_{I^*}$ :  $ILR(f_2)$  trained on infinite auto-labeled data

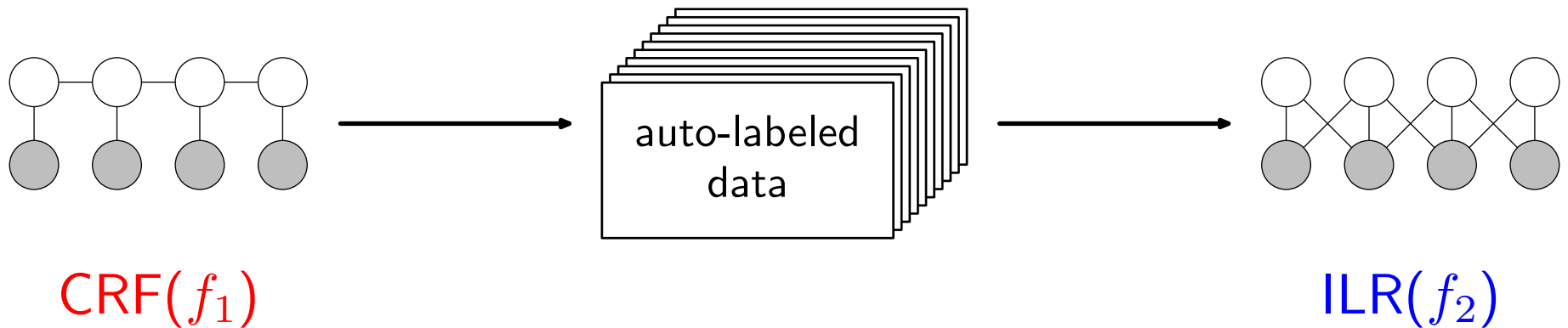
$$KL(p_C || p_I) = \underbrace{KL(p_C || p_{I^*})}_{\text{approx. error}} + \underbrace{(KL(p_C || p_I) - KL(p_C || p_{I^*}))}_{\text{estimation error}}$$

Estimation error:

$$\text{Expected value} = \frac{\# \text{ features}}{m} + o\left(\frac{1}{m}\right) \rightarrow 0 \text{ [Liang \& Jordan, 2008]}$$



# Approximation/estimation errors for ILR



$p_C$ :  $CRF(f_1)$  trained on labeled data

$p_I$ :  $ILR(f_2)$  trained on  $m$  auto-labeled examples

$p_{I^*}$ :  $ILR(f_2)$  trained on infinite auto-labeled data

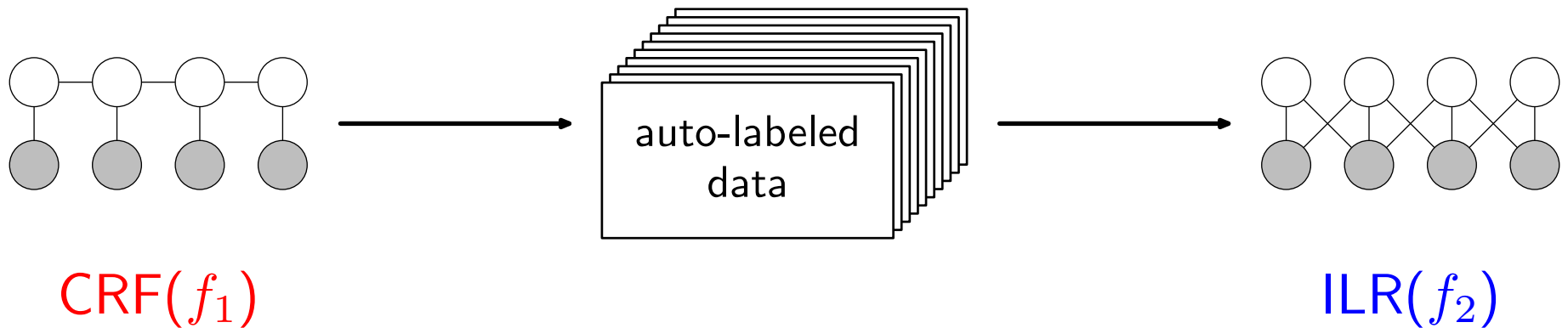
$$\text{KL} (p_C || p_I) = \underbrace{\text{KL} (p_C || p_{I^*})}_{\text{approx. error}} + \underbrace{(\text{KL} (p_C || p_I) - \text{KL} (p_C || p_{I^*}))}_{\text{estimation error}}$$

Estimation error:

Expected value =  $\frac{\# \text{ features}}{m} + o\left(\frac{1}{m}\right) \rightarrow 0$  [Liang & Jordan, 2008]

Structured compilation can eliminate this error

# Approximation/estimation errors for ILR



$p_C$ :  $CRF(f_1)$  trained on labeled data

$p_I$ :  $ILR(f_2)$  trained on  $m$  auto-labeled examples

$p_{I^*}$ :  $ILR(f_2)$  trained on infinite auto-labeled data

$$KL(p_C || p_I) = \underbrace{KL(p_C || p_{I^*})}_{\text{approx. error}} + \underbrace{(KL(p_C || p_I) - KL(p_C || p_{I^*}))}_{\text{estimation error}}$$

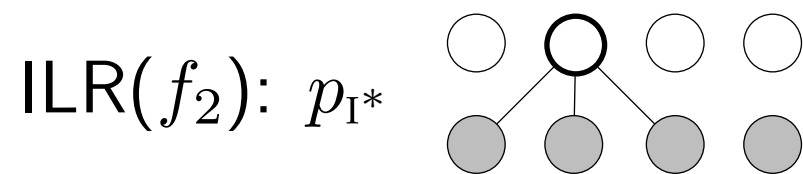
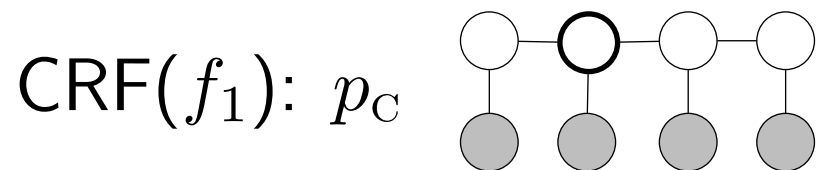
Estimation error:

Expected value =  $\frac{\# \text{ features}}{m} + o\left(\frac{1}{m}\right) \rightarrow 0$  [Liang & Jordan, 2008]

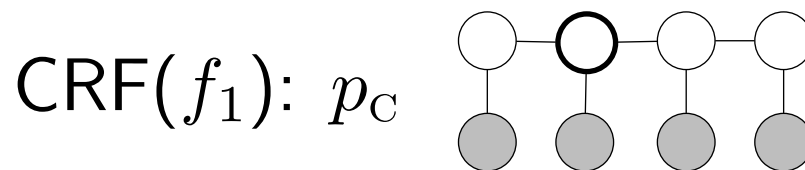
Structured compilation can eliminate this error

Approximation error: next...

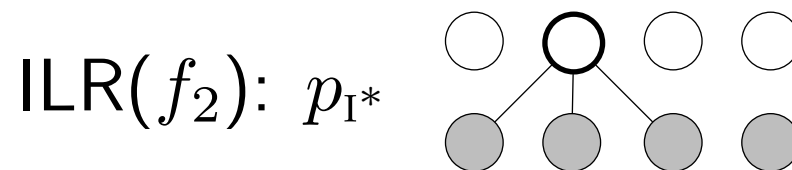
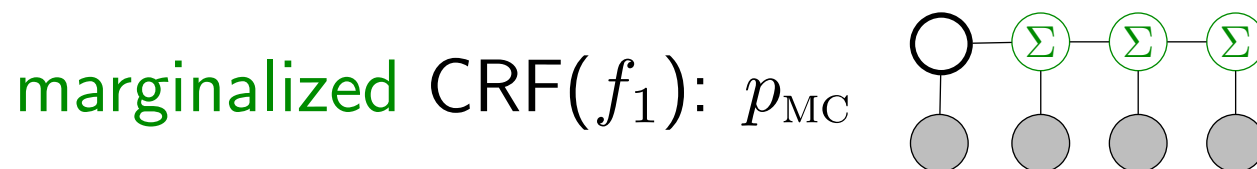
# Decomposition of approximation error



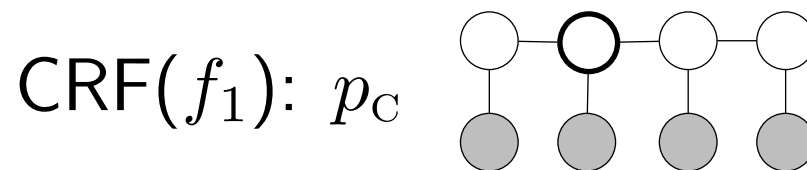
# Decomposition of approximation error



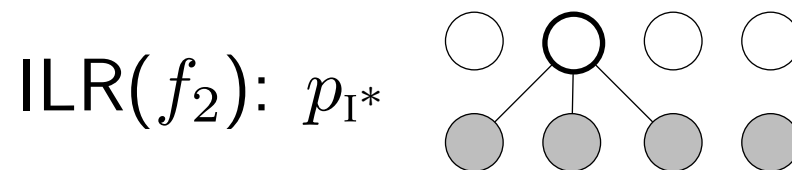
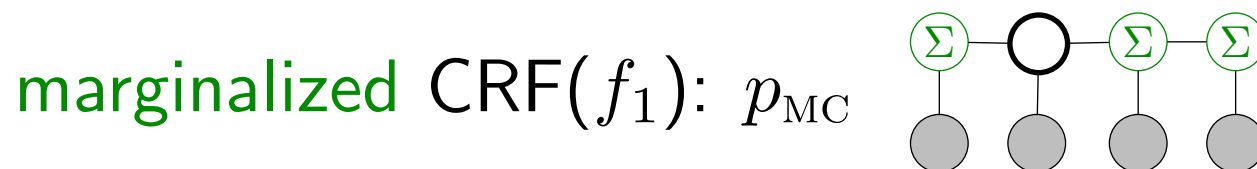
coherence



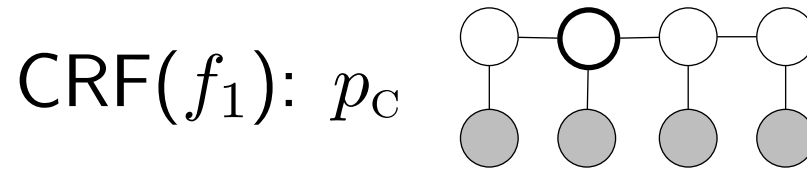
# Decomposition of approximation error



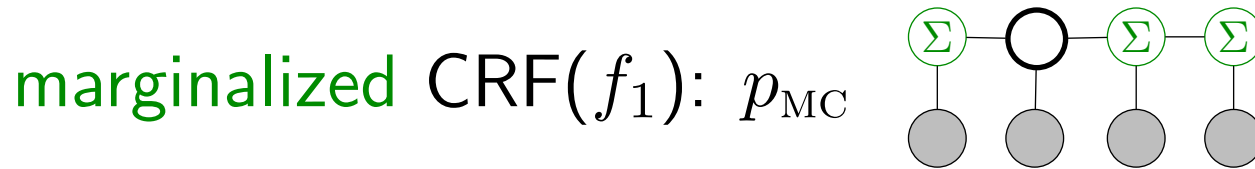
coherence



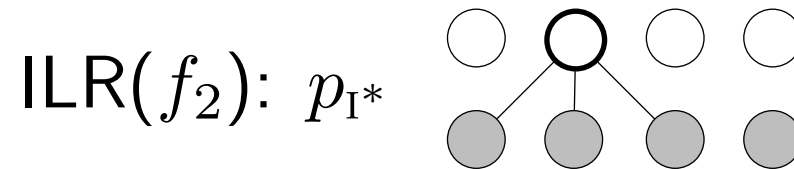
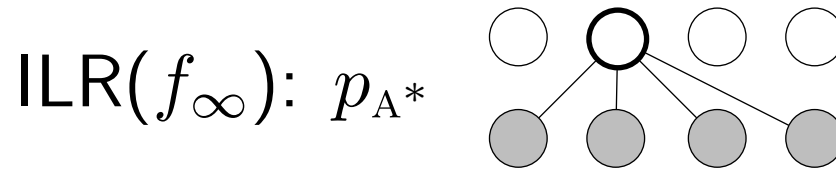
# Decomposition of approximation error



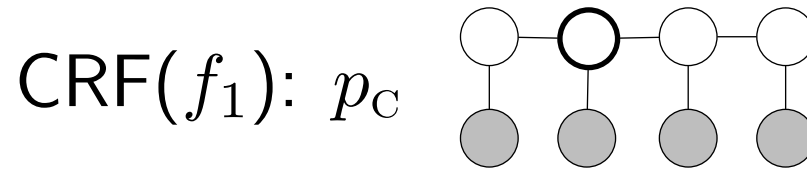
coherence



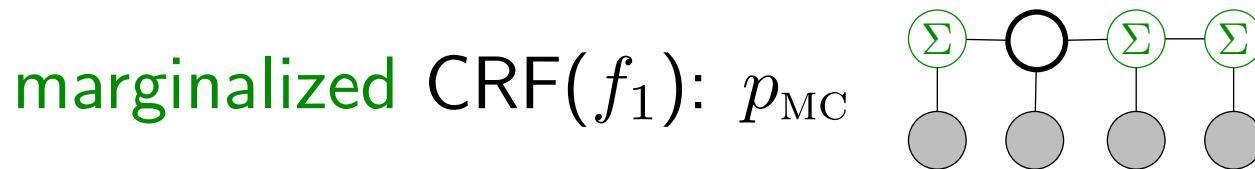
nonlinearities



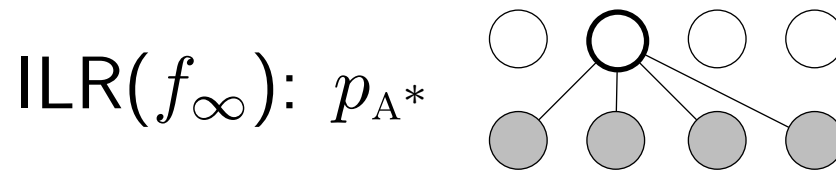
# Decomposition of approximation error



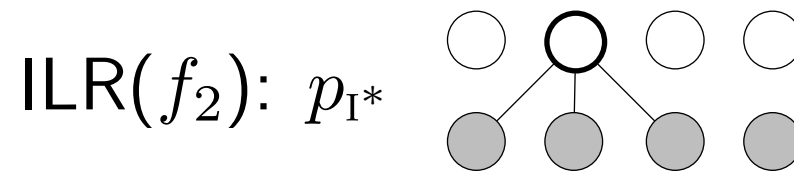
coherence



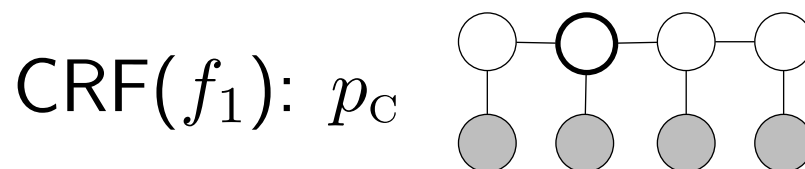
nonlinearities



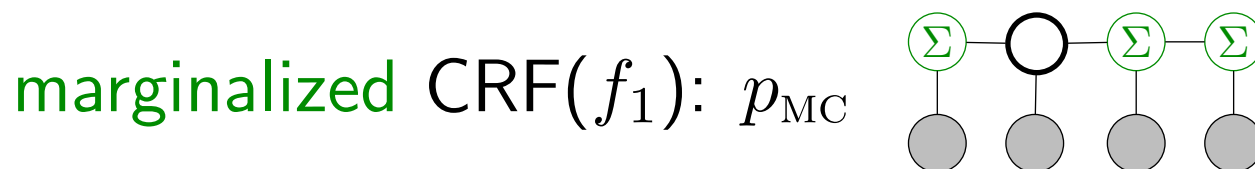
global information



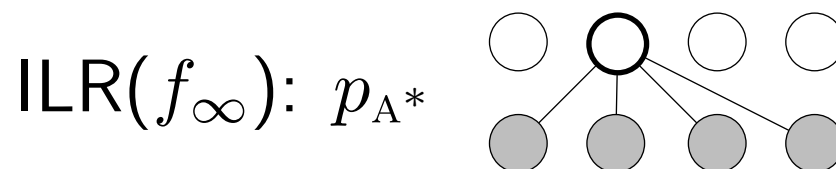
# Decomposition of approximation error



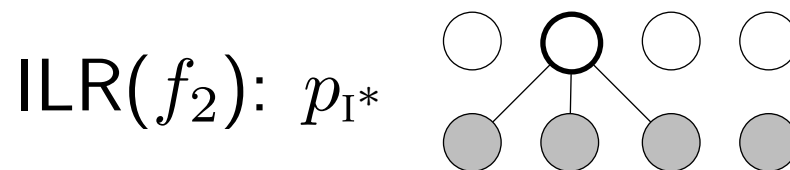
coherence



nonlinearities



global information

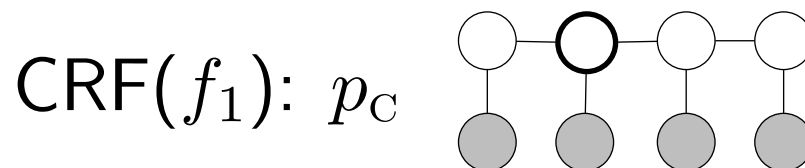


Theorem:

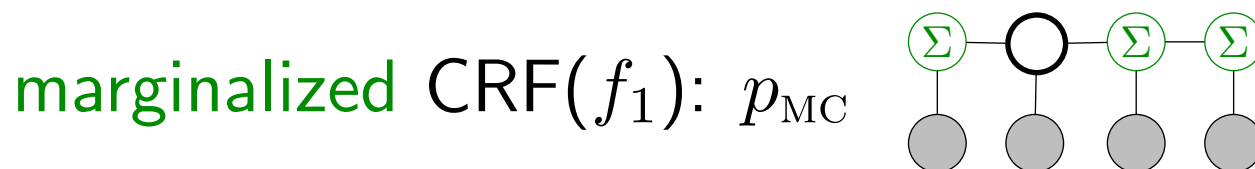
$$\text{KL}(p_C \parallel p_{I^*}) = \text{KL}(p_C \parallel p_{MC}) + \text{KL}(p_{MC} \parallel p_{A^*}) + \text{KL}(p_{A^*} \parallel p_{I^*})$$



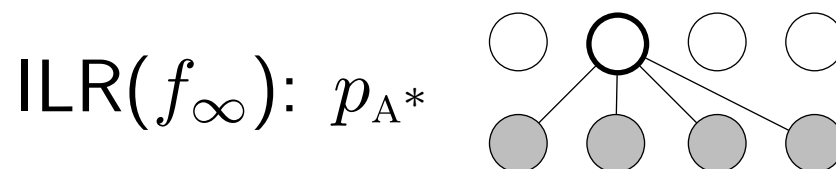
# Decomposition of approximation error



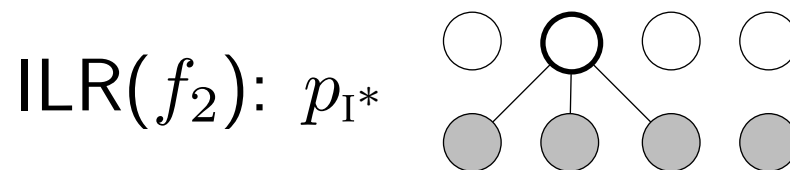
coherence



nonlinearities



global information



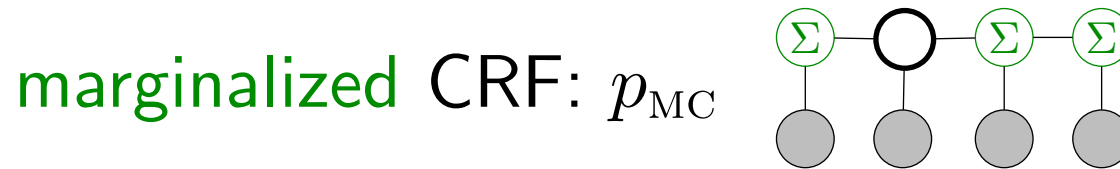
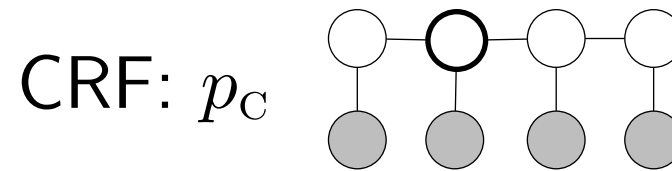
Theorem:

$$\text{KL}(p_C || p_{I^*}) = \text{KL}(p_C || p_{MC}) + \text{KL}(p_{MC} || p_{A^*}) + \text{KL}(p_{A^*} || p_{I^*})$$

Proof:

Generalized Pythagorean identity for KL-divergence

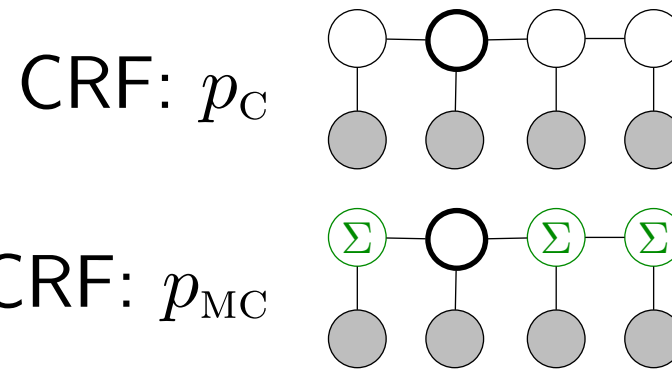
# Approximation error: coherence



**Coherence** =  $\text{KL}(p_C || p_{MC})$ :

importance of making joint predictions

# Approximation error: coherence



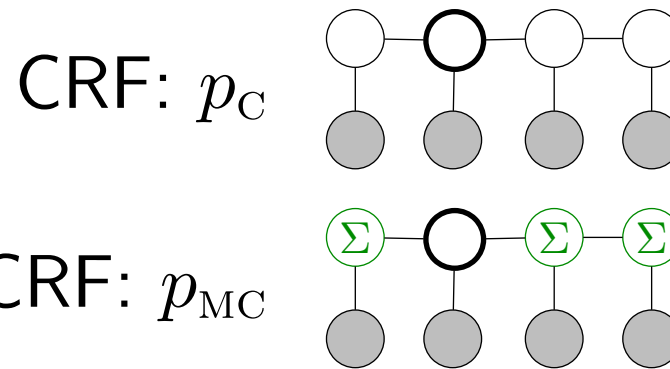
**Coherence** =  $\text{KL}(p_C || p_{MC})$ :

importance of making joint predictions

For a chain CRF:

coherence = sum of mutual information along the edges

# Approximation error: coherence



**Coherence** =  $\text{KL}(p_C || p_{MC})$ :

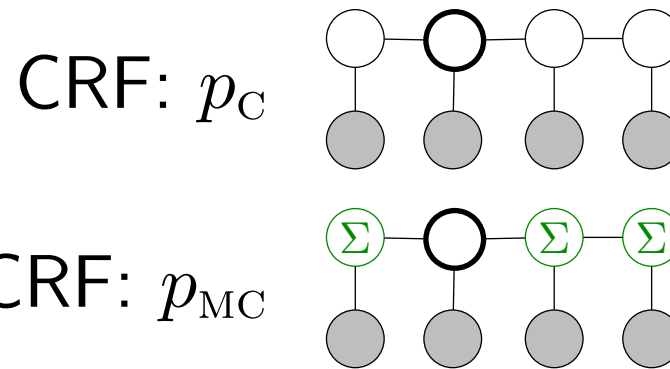
importance of making joint predictions

For a chain CRF:

coherence = sum of mutual information along the edges

	<b>POS</b>	<b>NER</b>
<b>Coherence</b>	0.003	0.009

# Approximation error: coherence



**Coherence** =  $\text{KL}(p_C || p_{MC})$ :

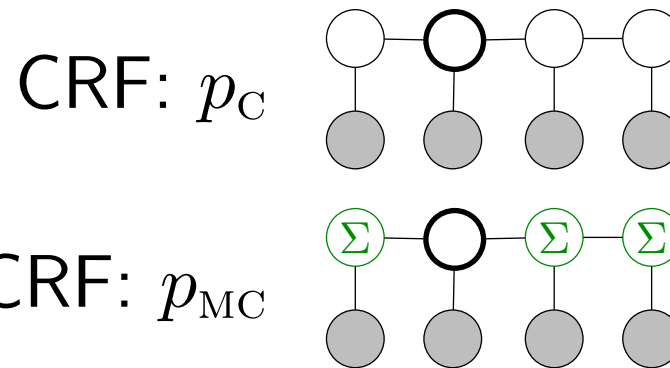
importance of making joint predictions

For a chain CRF:

coherence = sum of mutual information along the edges

	<b>POS</b>	<b>NER</b>
<b>Coherence</b>	0.003	0.009
<b>Change in accuracy</b>	95.0% $\Rightarrow$ 95.0%	76.3% $\Rightarrow$ 76.0%

# Approximation error: coherence



**Coherence** =  $\text{KL}(p_C || p_{MC})$ :

importance of making joint predictions

For a chain CRF:

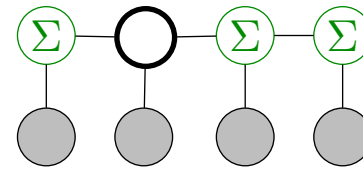
coherence = sum of mutual information along the edges

	<b>POS</b>	<b>NER</b>
<b>Coherence</b>	0.003	0.009
<b>Change in accuracy</b>	95.0% $\Rightarrow$ 95.0%	76.3% $\Rightarrow$ 76.0%

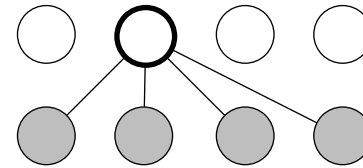
**Coherence is not a huge concern** (for these applications)

# Approximation error: nonlinearities

marginalized CRF:  $p_{MC}$



ILR( $f_\infty$ ):  $p_{A^*}$

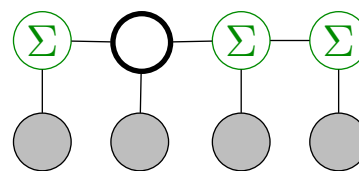


**Nonlinearities** =  $\text{KL}(p_{MC} || p_{A^*})$ :

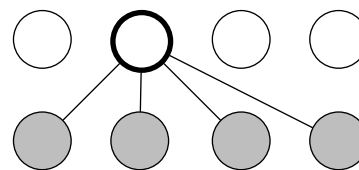
importance of combining features in a nonlinear way

# Approximation error: nonlinearities

marginalized CRF:  $p_{MC}$



ILR( $f_\infty$ ):  $p_{A^*}$



**Nonlinearities** =  $\text{KL}(p_{MC} || p_{A^*})$ :

importance of combining features in a nonlinear way

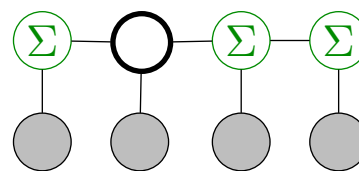
**NER experiment:**

Train a truncated CRF, so that both the truncated CRF (nonlinear) and the ILR (linear) use the **same features**

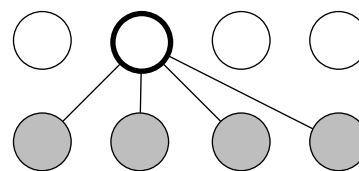


# Approximation error: nonlinearities

marginalized CRF:  $p_{MC}$



ILR( $f_\infty$ ):  $p_{A^*}$



**Nonlinearities** =  $\text{KL}(p_{MC} || p_{A^*})$ :

importance of combining features in a nonlinear way

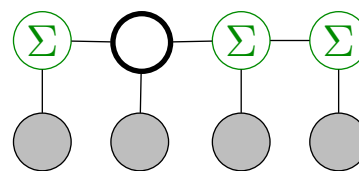
**NER experiment:**

Train a truncated CRF, so that both the truncated CRF (nonlinear) and the ILR (linear) use the **same features**

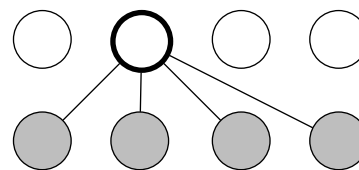
	<b>Truncated CRF</b>	<b>ILR(<math>f_2</math>)</b>
<b>Accuracy</b>	76.0%	72.7%

# Approximation error: nonlinearities

marginalized CRF:  $p_{MC}$



ILR( $f_\infty$ ):  $p_{A^*}$



**Nonlinearities** =  $\text{KL}(p_{MC} || p_{A^*})$ :

importance of combining features in a nonlinear way

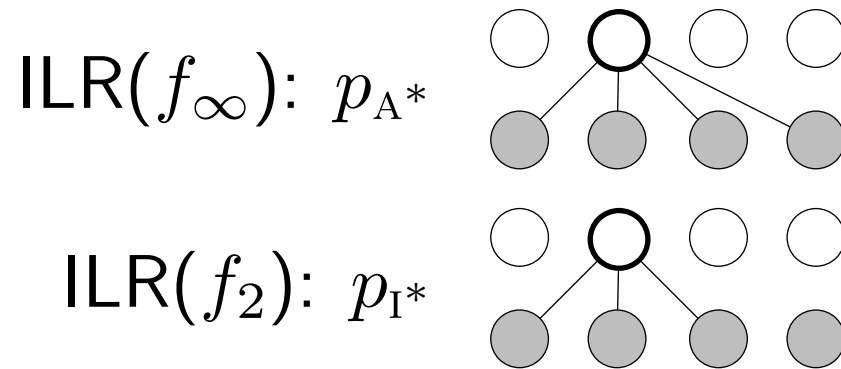
**NER experiment:**

Train a truncated CRF, so that both the truncated CRF (nonlinear) and the ILR (linear) use the **same features**

	<b>Truncated CRF</b>	<b>ILR(<math>f_2</math>)</b>
<b>Accuracy</b>	76.0%	72.7%

**Nonlinearities play an important role**

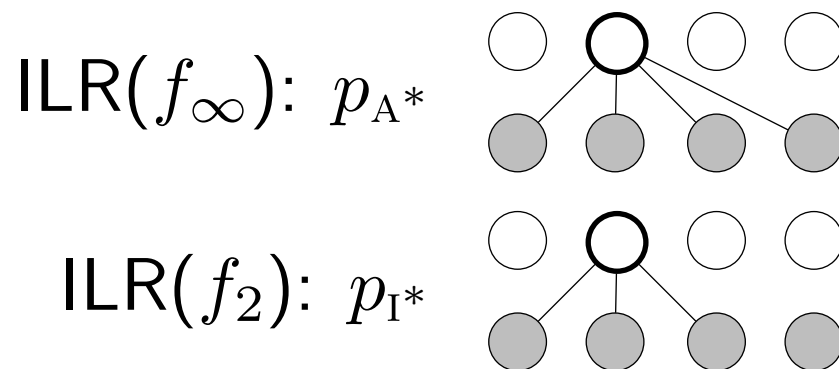
# Approximation error: global information



**Global information** =  $\text{KL}(p_{A^*} || p_{I^*})$ :

importance of using features on distant parts of the input

# Approximation error: global information



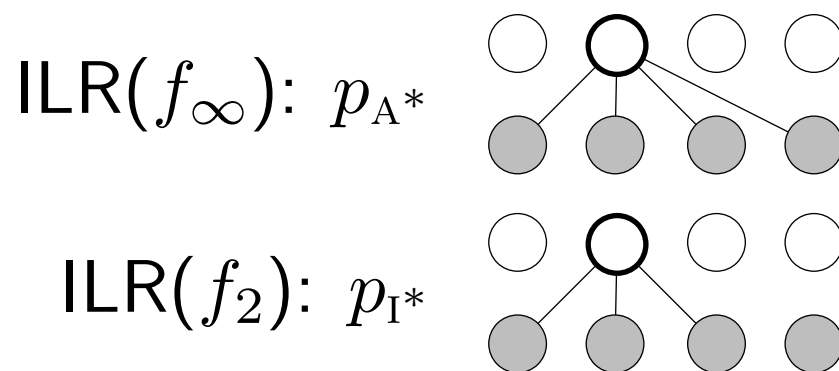
**Global information** =  $\text{KL}(p_{A^*} || p_{I^*})$ :

importance of using features on distant parts of the input

**NER experiment:**

Compare truncated CRF with marginalized CRF (they differ only in the features used)

# Approximation error: global information



**Global information** =  $\text{KL}(p_{A^*} || p_{I^*})$ :

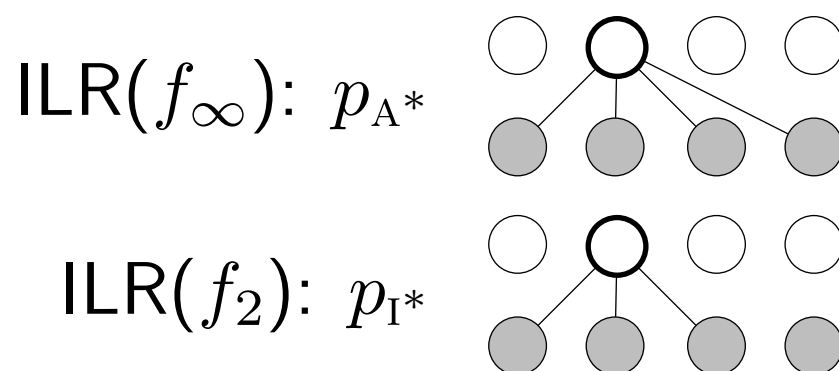
importance of using features on distant parts of the input

**NER experiment:**

Compare truncated CRF with marginalized CRF (they differ only in the features used)

	<b>Marginalized CRF</b>	<b>Truncated CRF</b>
<b>Accuracy</b>	76.0%	76.0%

# Approximation error: global information



**Global information** =  $\text{KL}(p_{A^*} || p_{I^*})$ :

importance of using features on distant parts of the input

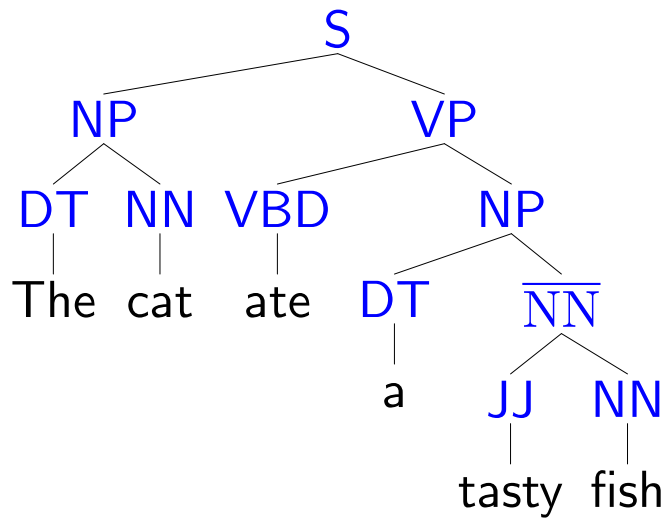
**NER experiment:**

Compare truncated CRF with marginalized CRF (they differ only in the features used)

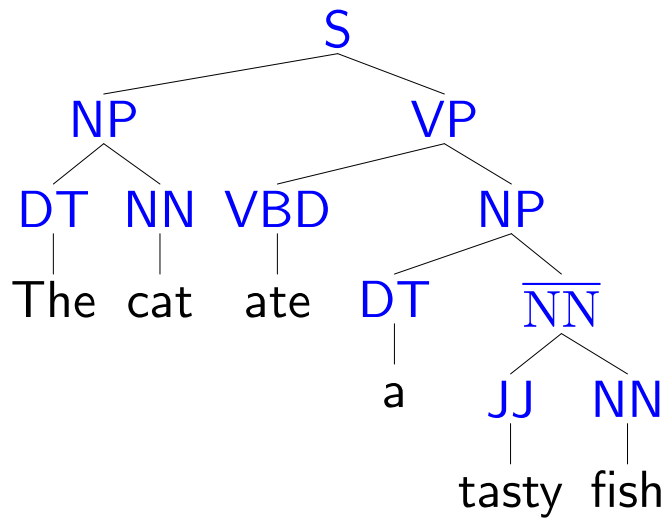
	<b>Marginalized CRF</b>	<b>Truncated CRF</b>
<b>Accuracy</b>	76.0%	76.0%

**Distant information is not essential** (for these applications)

# Structure compilation for parsing



# Structure compilation for parsing



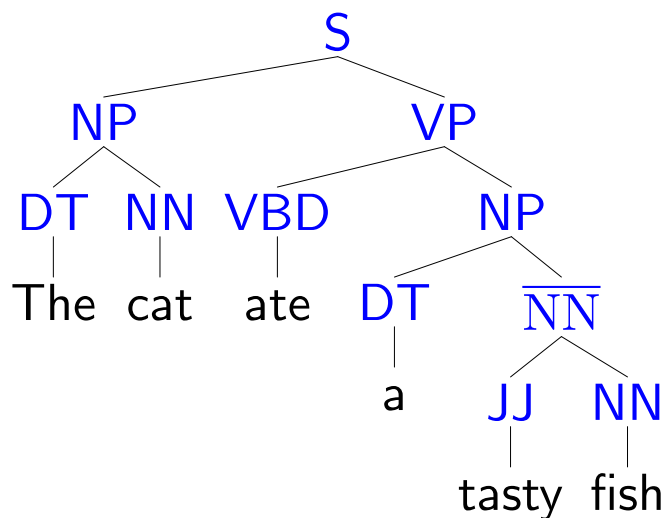
Sentence length:  $\ell$

Number of grammar symbols:  $K$

Number of grammar rules:  $G \gg \ell, K$



# Structure compilation for parsing



Sentence length:  $\ell$

Number of grammar symbols:  $K$

Number of grammar rules:  $G \gg \ell, K$

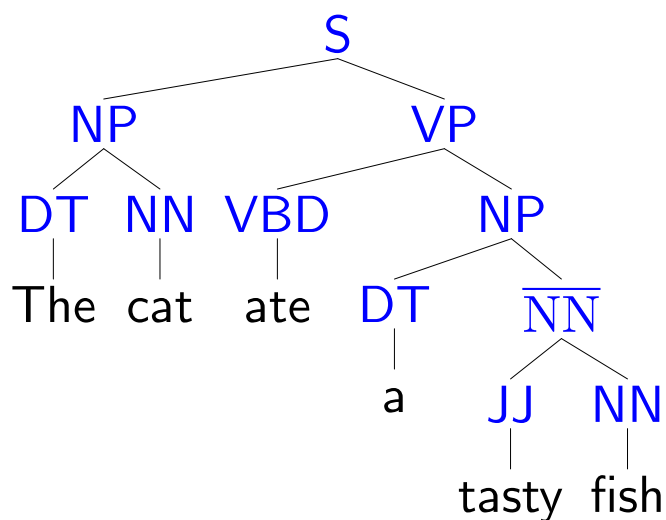
**Parse time/sentence**

$$O(\ell^3 G)$$

**Structured model:**

Standard dynamic program for context-free grammars

# Structure compilation for parsing



Sentence length:  $\ell$

Number of grammar symbols:  $K$

Number of grammar rules:  $G \gg \ell, K$

## Parse time/sentence

### Structured model:

$$O(\ell^3 G)$$

Standard dynamic program for context-free grammars

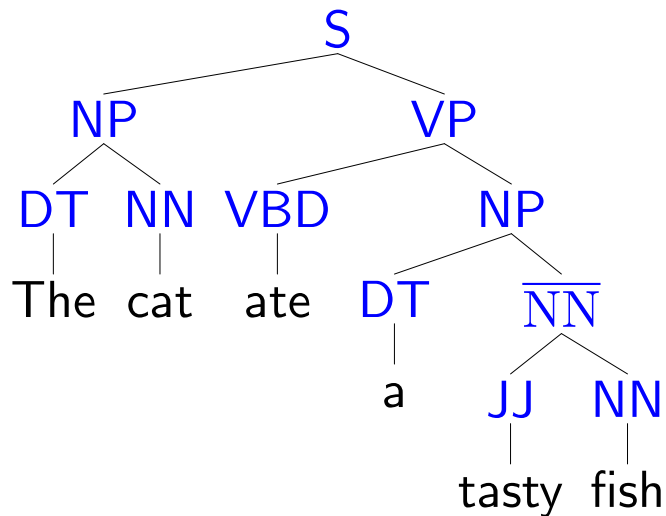
### Independent model:

$$O(\ell^3 + K\ell^2)$$

For each of  $O(\ell^2)$  spans:

Make a soft prediction of whether it's a constituent  
(features: words/tags/prefixes/suffixes on entire span)

# Structure compilation for parsing



Sentence length:  $\ell$

Number of grammar symbols:  $K$

Number of grammar rules:  $G \gg \ell, K$

## Parse time/sentence

### Structured model:

$$O(\ell^3 G)$$

Standard dynamic program for context-free grammars

### Independent model:

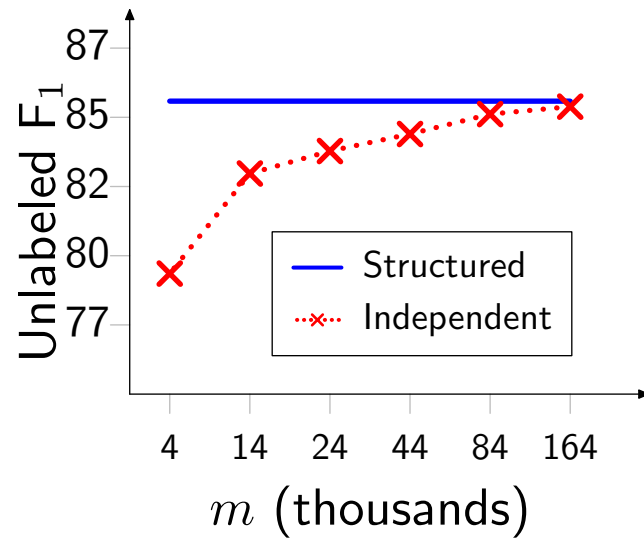
$$O(\ell^3 + K\ell^2)$$

For each of  $O(\ell^2)$  spans:

Make a soft prediction of whether it's a constituent  
(features: words/tags/prefixes/suffixes on entire span)

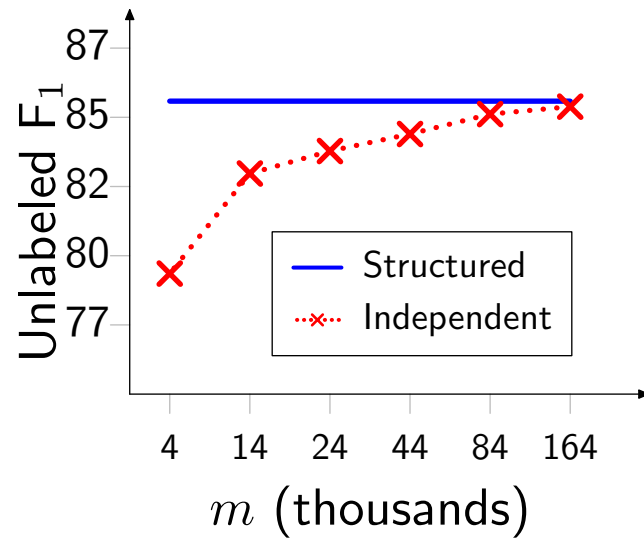
Run a dynamic program to choose the best tree

# Parsing results

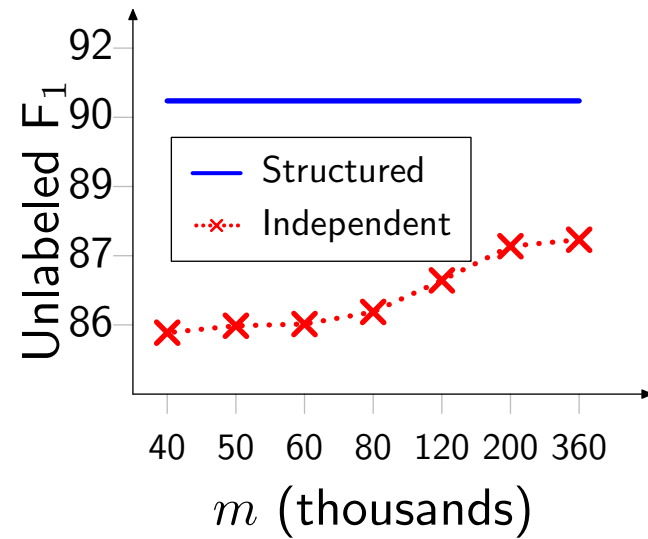


4K labeled sentences

# Parsing results

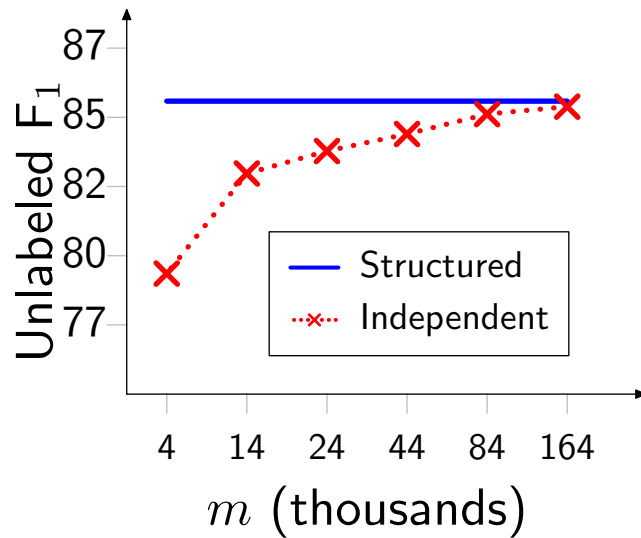


4K labeled sentences

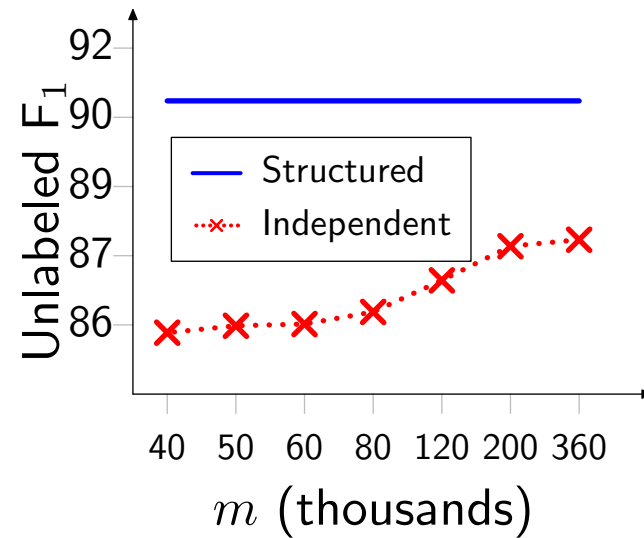


40K labeled sentences

# Parsing results



4K labeled sentences



40K labeled sentences

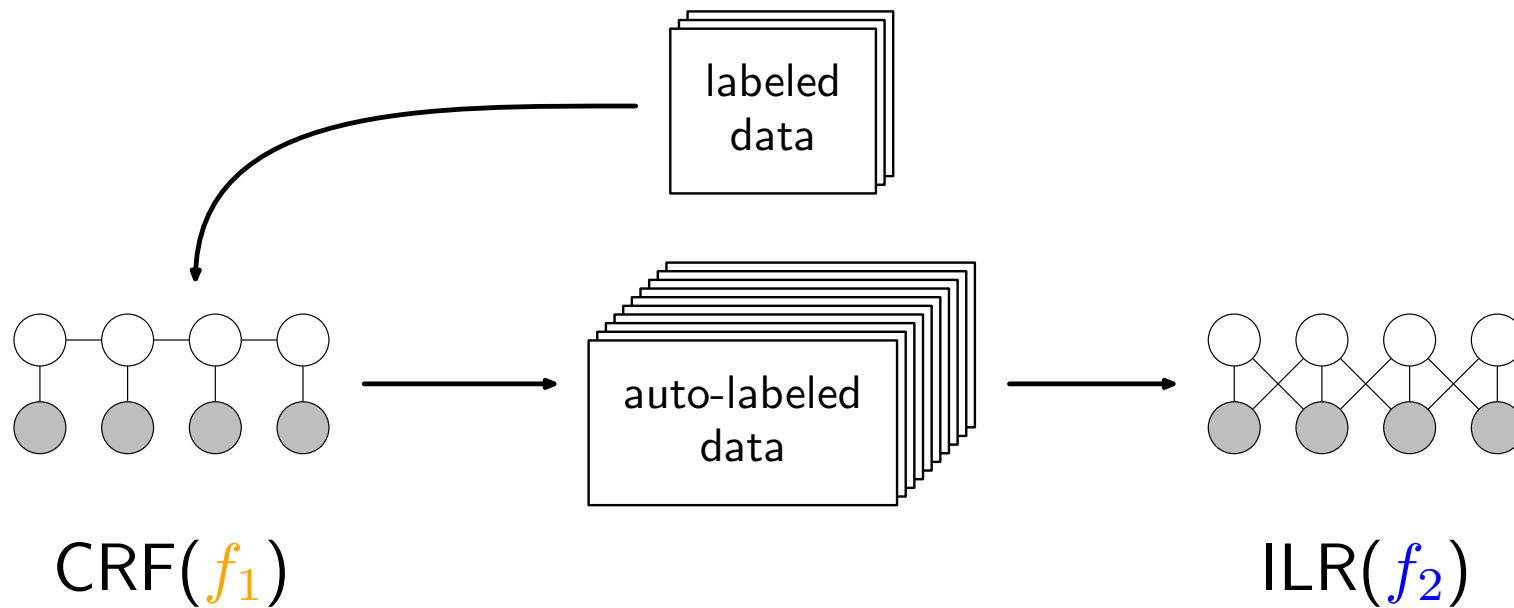
- Structure is important in parsing
- Need richer features or nonlinearities for the independent model to catch up

# Summary of structure compilation

Motivation: want fast CRF-level accuracy at test time

# Summary of structure compilation

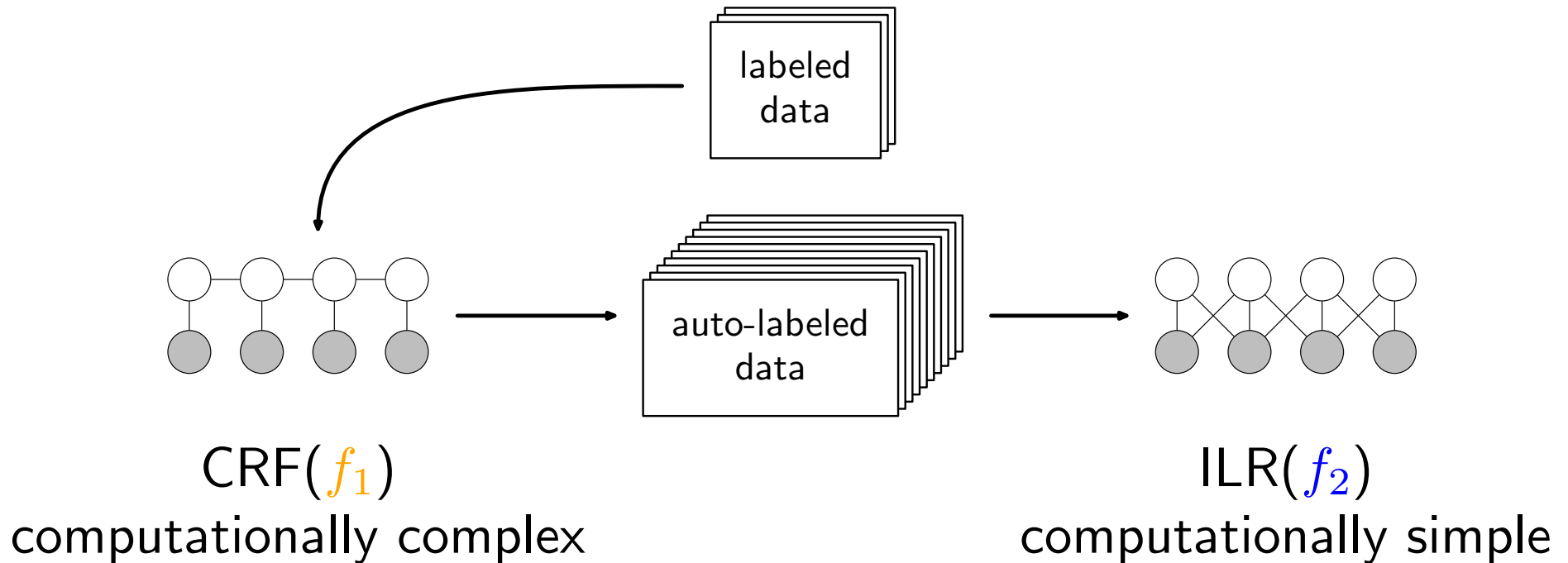
Motivation: want fast CRF-level accuracy at test time





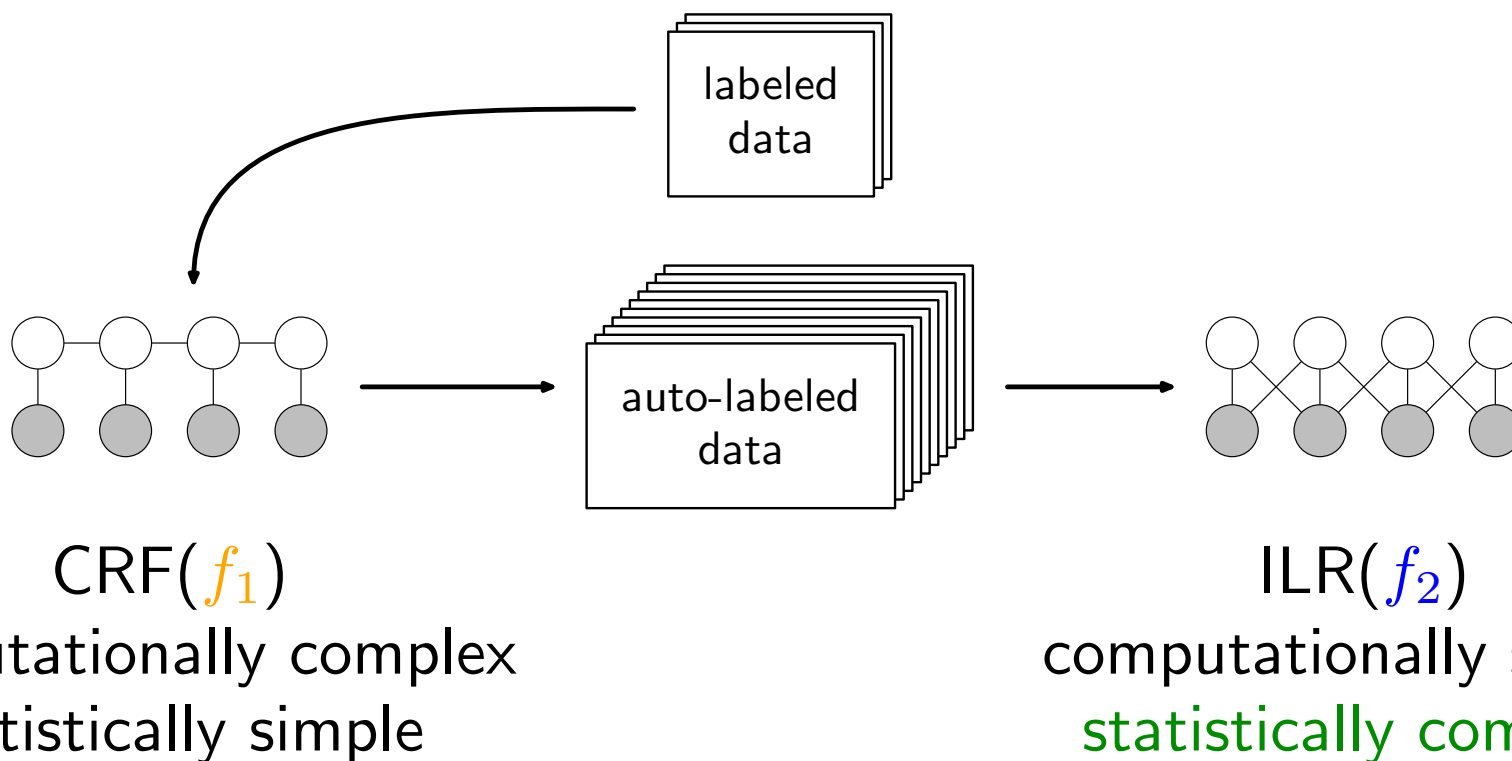
# Summary of structure compilation

Motivation: want fast CRF-level accuracy at test time



# Summary of structure compilation

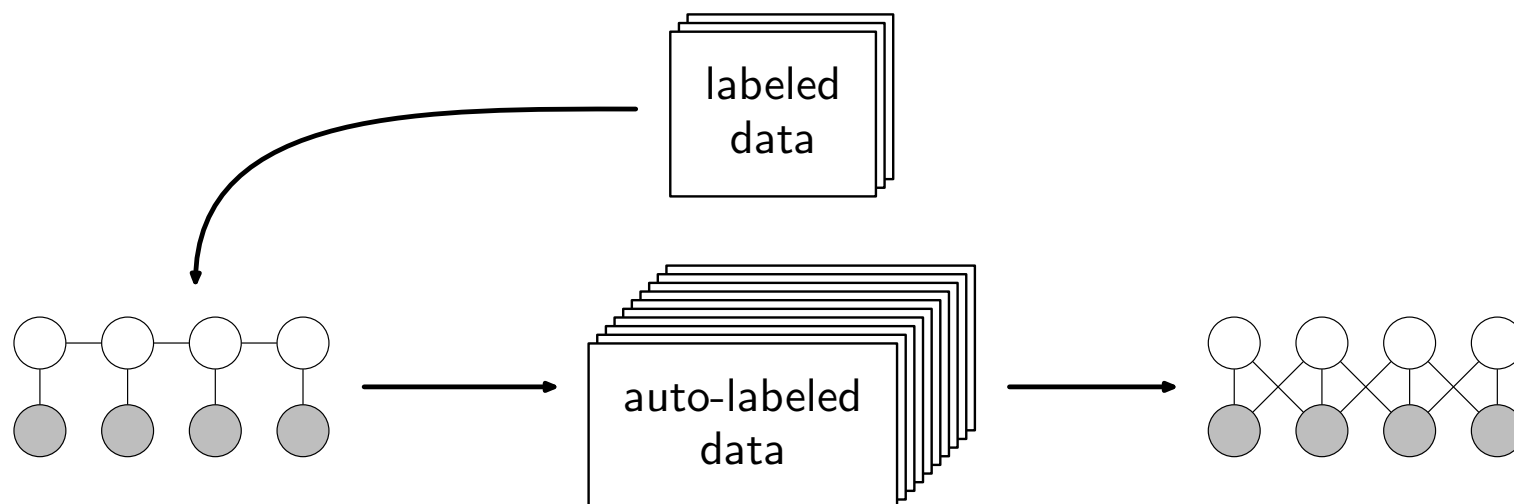
Motivation: want fast CRF-level accuracy at test time



Estimation error: structure compilation can easily drive it to 0

# Summary of structure compilation

Motivation: want fast CRF-level accuracy at test time



CRF( $f_1$ )

computationally complex  
statistically simple  
very expressive

ILR( $f_2$ )

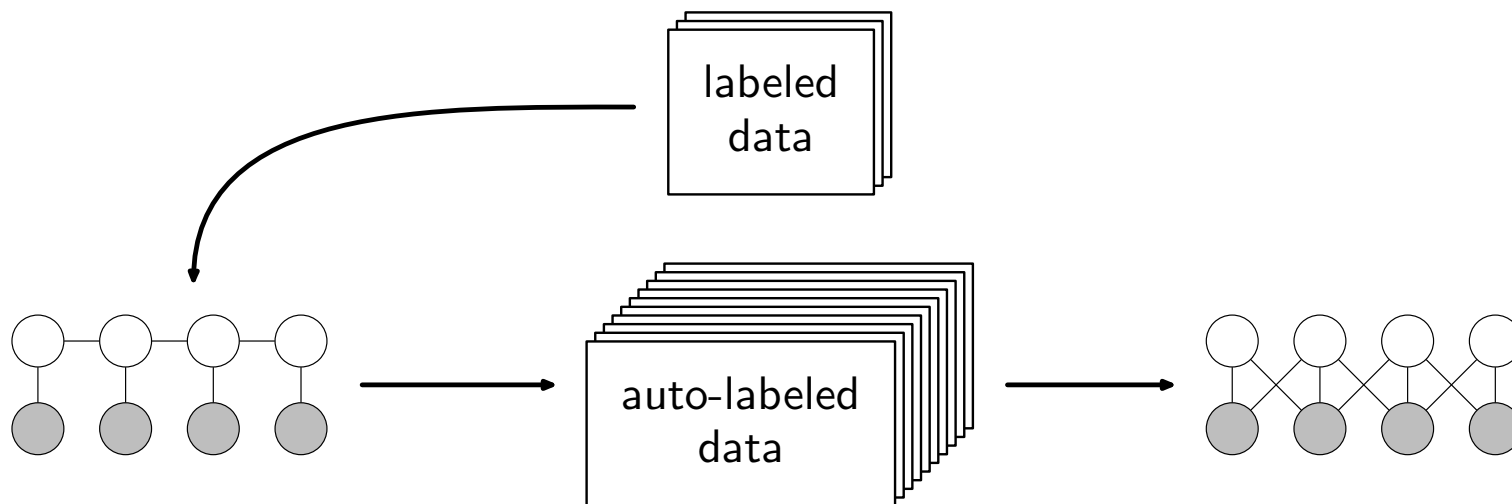
computationally simple  
statistically complex  
not as expressive

**Estimation error:** structure compilation can easily drive it to 0

**Approximation error:** advantages of CRF over ILR

# Summary of structure compilation

Motivation: want fast CRF-level accuracy at test time



CRF( $f_1$ )

computationally complex  
statistically simple  
very expressive

ILR( $f_2$ )

computationally simple  
statistically complex  
not as expressive

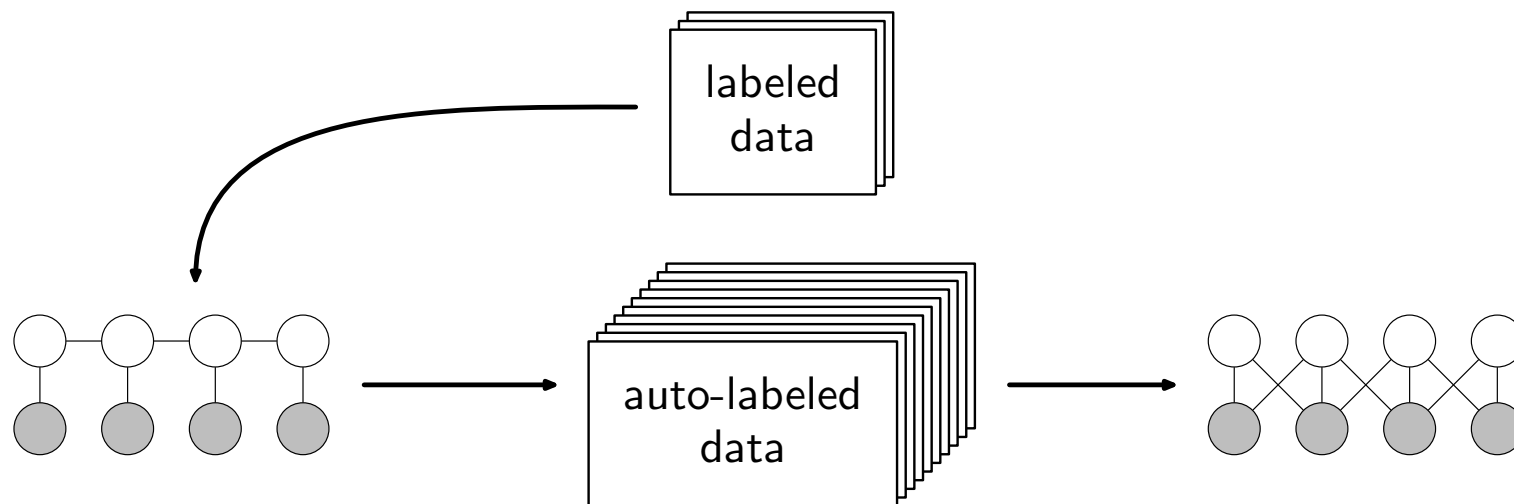
**Estimation error:** structure compilation can easily drive it to 0

**Approximation error:** advantages of CRF over ILR

- ILR needs rich features to compensate

# Summary of structure compilation

Motivation: want fast CRF-level accuracy at test time



CRF( $f_1$ )

computationally complex  
statistically simple  
very expressive

ILR( $f_2$ )

computationally simple  
statistically complex  
not as expressive

**Estimation error:** structure compilation can easily drive it to 0

**Approximation error:** advantages of CRF over ILR

- ILR needs rich features to compensate
- CRF's nonlinearities are important