
Domain Adaptation

Philipp Koehn

30 November 2017



"Domain"



- Corpora differ
 - topic (politics, news, medicine, ...)
 - style (formal, informal)
 - modality (written, transcribed speech)
 - register (level of politeness)
- Covered on the catch-all term "domain"
- Domain := one source for a parallel corpus

"Domain"



- Domain matters for word choice
 - *bat* in baseball domain vs. *bat* in animal domain
 - *interest* in financial domain vs. *interest* in arts
- Style matters, too
 - translate greeting into *What's up?* vs. *Ladies and Gentlemen!*
 - use of informal *Du* vs. formal *Sie* in German
- Distinctions often only visible in full document / full corpus

Various Data Sources

- Available parallel corpora on OPUS web site (Italian–English)

corpus	doc's	sent's	it tokens	en tokens	XCES/XML	raw	TMX	Moses
OpenSubtitles2018	48746	37.8M	304.8M	284.5M	[xces en it]	[en it]	[tmx]	[mooses]
EUbookshop	9028	6.6M	268.7M	258.8M	[xces en it]	[en it]	[tmx]	[mooses]
OpenSubtitles2016	35929	28.7M	230.3M	214.9M	[xces en it]	[en it]	[tmx]	[mooses]
DGT	26880	3.2M	72.9M	64.0M	[xces en it]	[en it]	[tmx]	[mooses]
Europarl	9461	2.0M	59.9M	58.9M	[xces en it]	[en it]	[tmx]	[mooses]
JRC-Acquis	12042	0.8M	34.1M	34.5M	[xces en it]	[en it]	[tmx]	[mooses]
Wikipedia	3	1.0M	26.5M	22.2M	[xces en it]	[en it]	[tmx]	[mooses]
EMEA	1920	1.1M	12.0M	13.9M	[xces en it]	[en it]	[tmx]	[mooses]
ECB	1	0.2M	5.5M	5.8M	[xces en it]	[en it]	[tmx]	[mooses]
GNOME	1905	0.7M	3.8M	3.4M	[xces en it]	[en it]	[tmx]	[mooses]
TED2013	1	0.2M	3.2M	2.7M	[xces en it]	[en it]	[tmx]	[mooses]
Tanzil	15	0.1M	2.8M	2.4M	[xces en it]	[en it]	[tmx]	[mooses]
Tatoeba	1	0.1M	3.6M	1.3M	[xces en it]	[en it]	[tmx]	[mooses]
KDE4	1957	0.3M	2.2M	2.3M	[xces en it]	[en it]	[tmx]	[mooses]
GlobalVoices	3220	81.3k	2.1M	2.0M	[xces en it]	[en it]	[tmx]	[mooses]
News-Commentary11	1423	45.9k	1.3M	1.0M	[xces en it]	[en it]	[tmx]	[mooses]
Books	8	33.1k	0.9M	0.8M	[xces en it]	[en it]	[tmx]	[mooses]
Ubuntu	452	0.1M	0.8M	0.6M	[xces en it]	[en it]	[tmx]	[mooses]
News-Commentary	1	18.6k	0.5M	0.5M	[xces en it]	[en it]	[tmx]	[mooses]
PHP	3270	36.8k	0.5M	0.2M	[xces en it]	[en it]	[tmx]	[mooses]
EUconst	47	10.2k	0.2M	0.2M	[xces en it]	[en it]	[tmx]	[mooses]
OpenSubtitles	22	19.1k	0.2M	0.1M	[xces en it]	[en it]	[tmx]	[mooses]
total	156332	83.1M	1.0G	975.1M	83.1M		63.4M	77.4M

Domain Examples



EMEA Abilify is a medicine containing the active substance aripiprazole.

It is available as 5 mg, 10 mg, 15 mg and 30 mg tablets, as 10 mg, 15 mg and 30 mg orodispersible tablets (tablets that dissolve in the mouth), as an oral solution (1 mg/ml) and as a solution for injection (7.5 mg/ml).

Software Localization Default GNOME Theme

OK

People

Pictures

Plan

Sound

Literature There was a slight noise behind her and she turned just in time to seize a small boy by the slack of his roundabout and arrest his flight.

Law Corrigendum to the Interim Agreement with a view to an Economic Partnership Agreement between the European Community and its Member States, of the one part, and the Central Africa Party, of the other part.

Domain Examples



PHP If you would like to start a new translation, or help in a translation project, please read <http://cvs.php.net/co.php/phpdoc/howto/howto.html.tar.gz>.

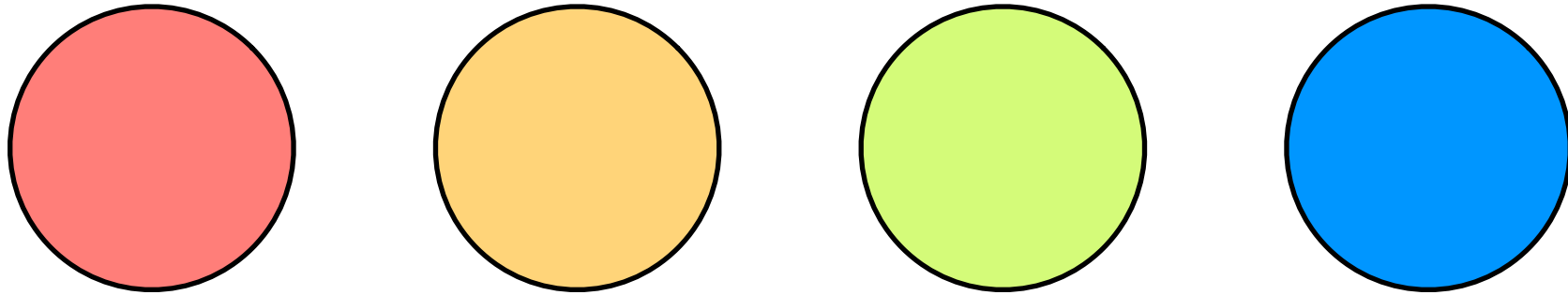
Religion This is The Book free of doubt and involution, a guidance for those who preserve themselves from evil and follow the straight path.

News The Facebook page of a leading Iranian leading cartoonist, Mana Nayestani, was hacked on Tuesday, 11 September 2012, by pro-regime hackers who call themselves "Soldiers of Islam".

Movie subtitles We're taking you to Washington, D.C.
Do you know where the prisoner was transported to?
Uh, Washington.
Okay.

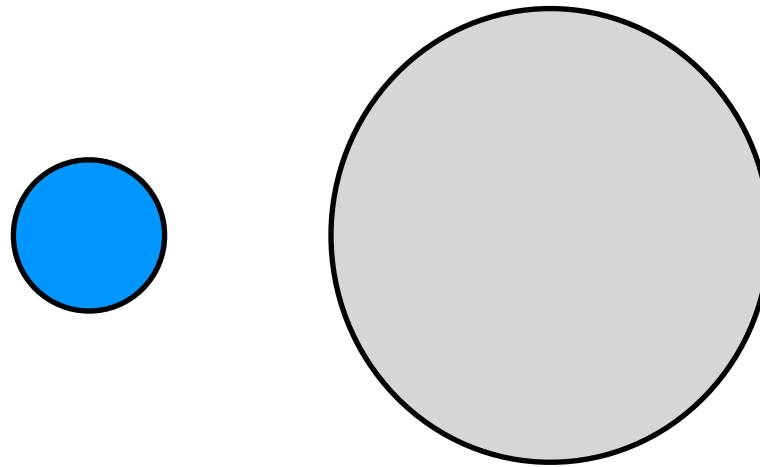
Twitter Thank u @Starbucks & @Spotify for celebrating artists who #GiveGood with a donation to @BTWFoundation, and to great organizations by @Metallica and @ChanceTheRapper! Limited edition cards available now at Starbucks!

Multi-Domain Scenario



- Machine translation systems work best when optimized for one domain
- Separate data by domain
- Build special system for each domain
- Translate each sentence with matching system

In/Out Domain Scenario



- Optimize system for just one domain
- Available data
 - small amounts of in-domain data
 - large amounts of out-of-domain data
- Need to balance both data sources

Why Use Out-of-Domain Data?



- In-domain data much more valuable
- But: gaps
 - word-to-be-translated may not occur
 - word-to-be-translated may not occur with the correct translation
- Motivation
 - out-of-domain data may fill these gaps
 - but be careful not to drown out in-domain data

S^4 Taxonomy of Adaptation Effects



9

[Carpuat, Daume, Fraser, Quirk, 2012]

- **Seen:** Never seen this word before

News to medical: diabetes mellitus

- **Sense:** Never seen this word used in this way

News to technical: monitor

- **Score:** The wrong output is scored higher

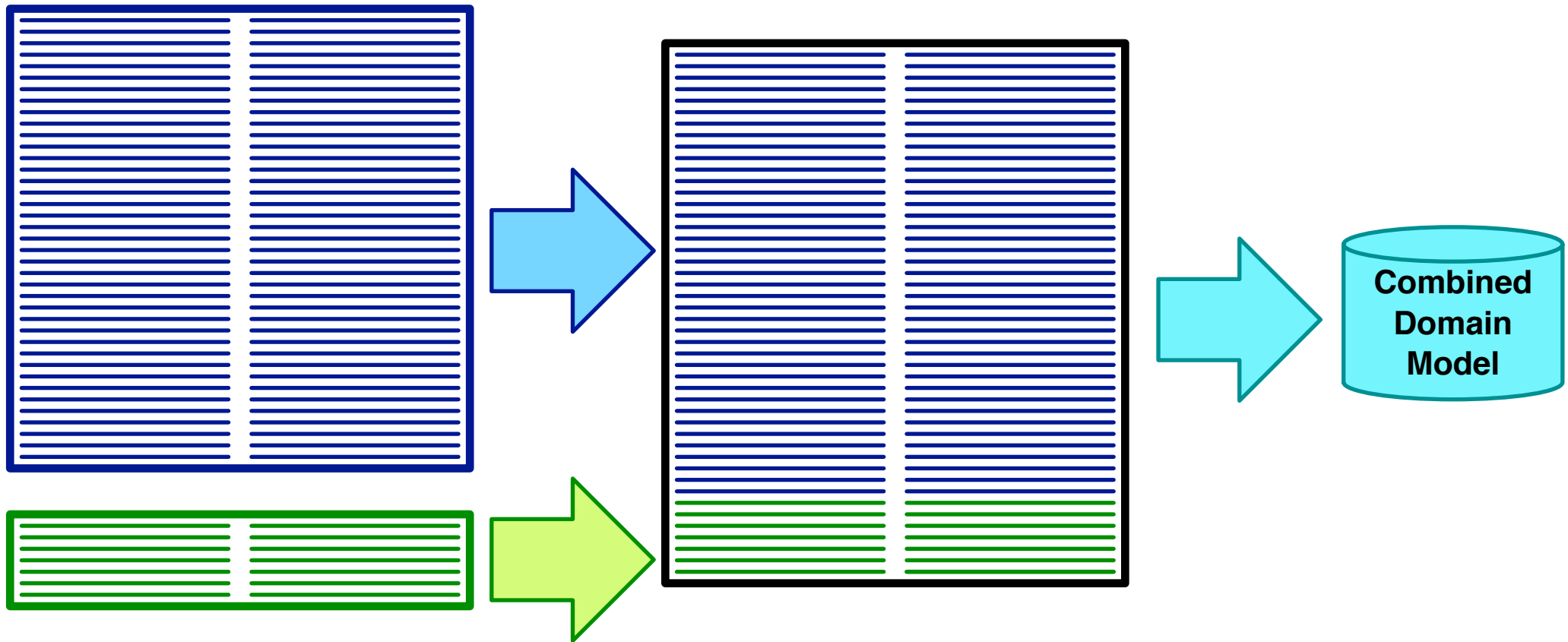
News to medical: manifest

- **Search:** Decoding/search erred



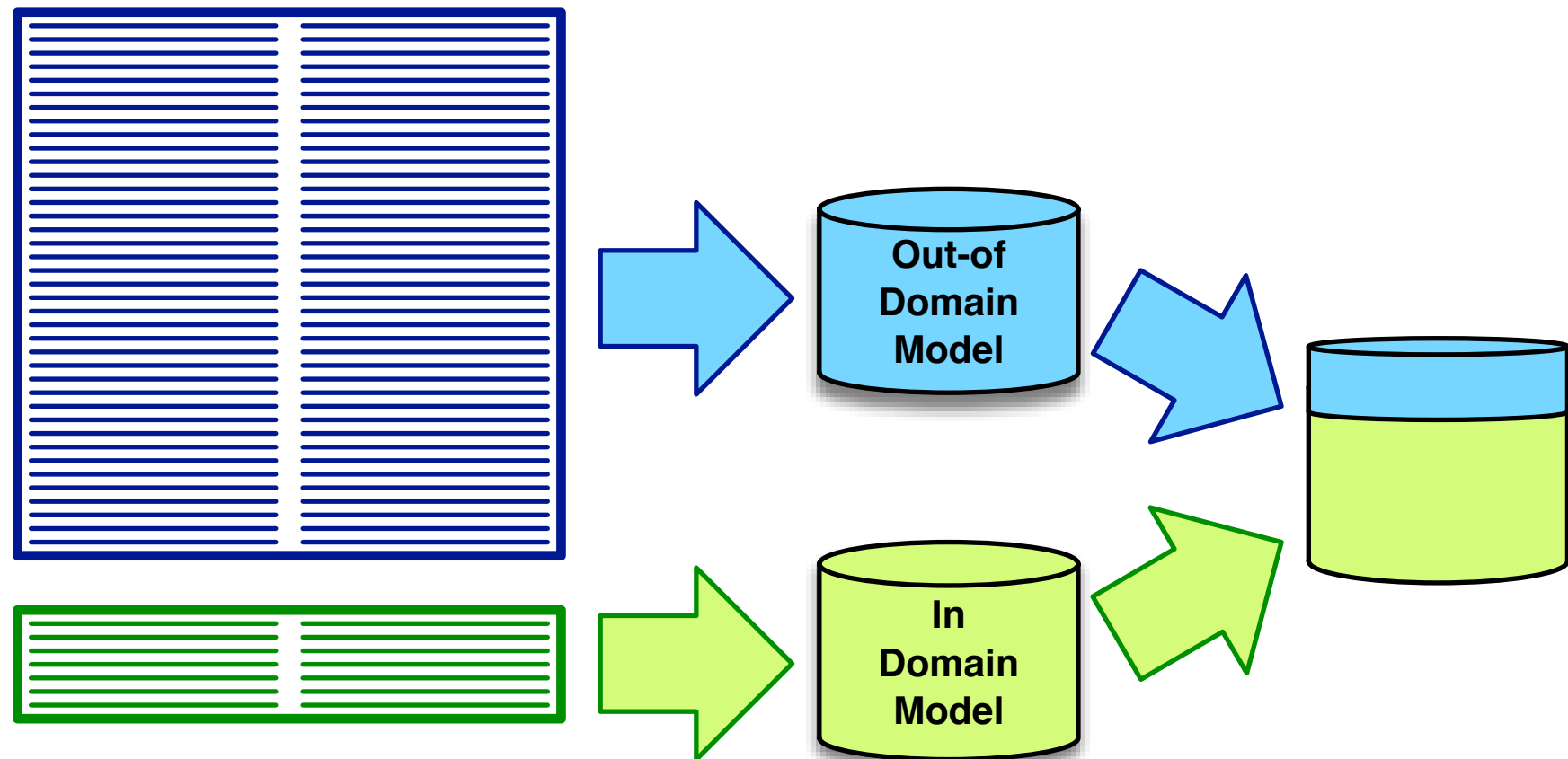
mixture models

Combining Data



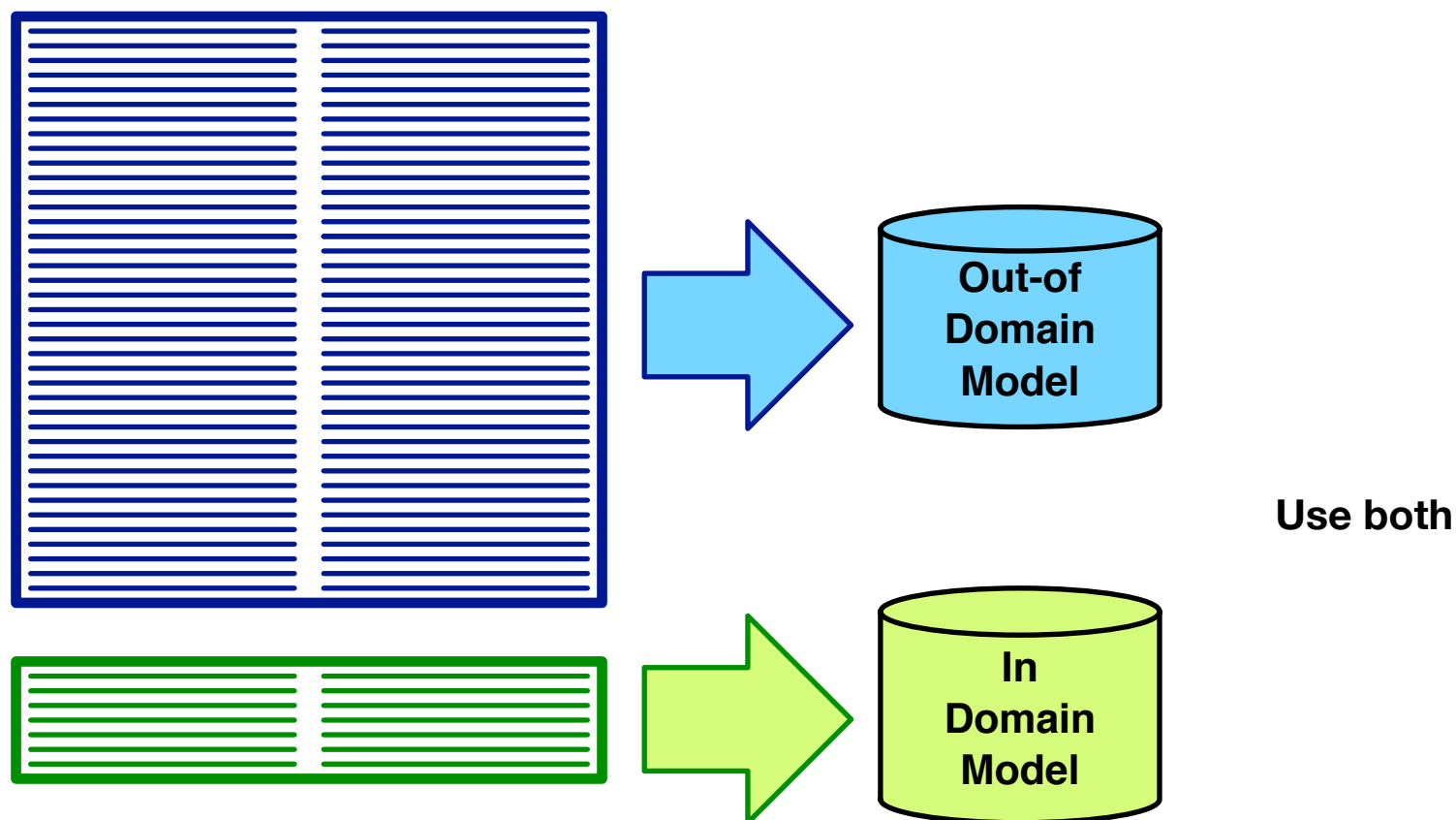
- Too biased towards out of domain data
- May flag translation options with indicator feature functions

Interpolate Models



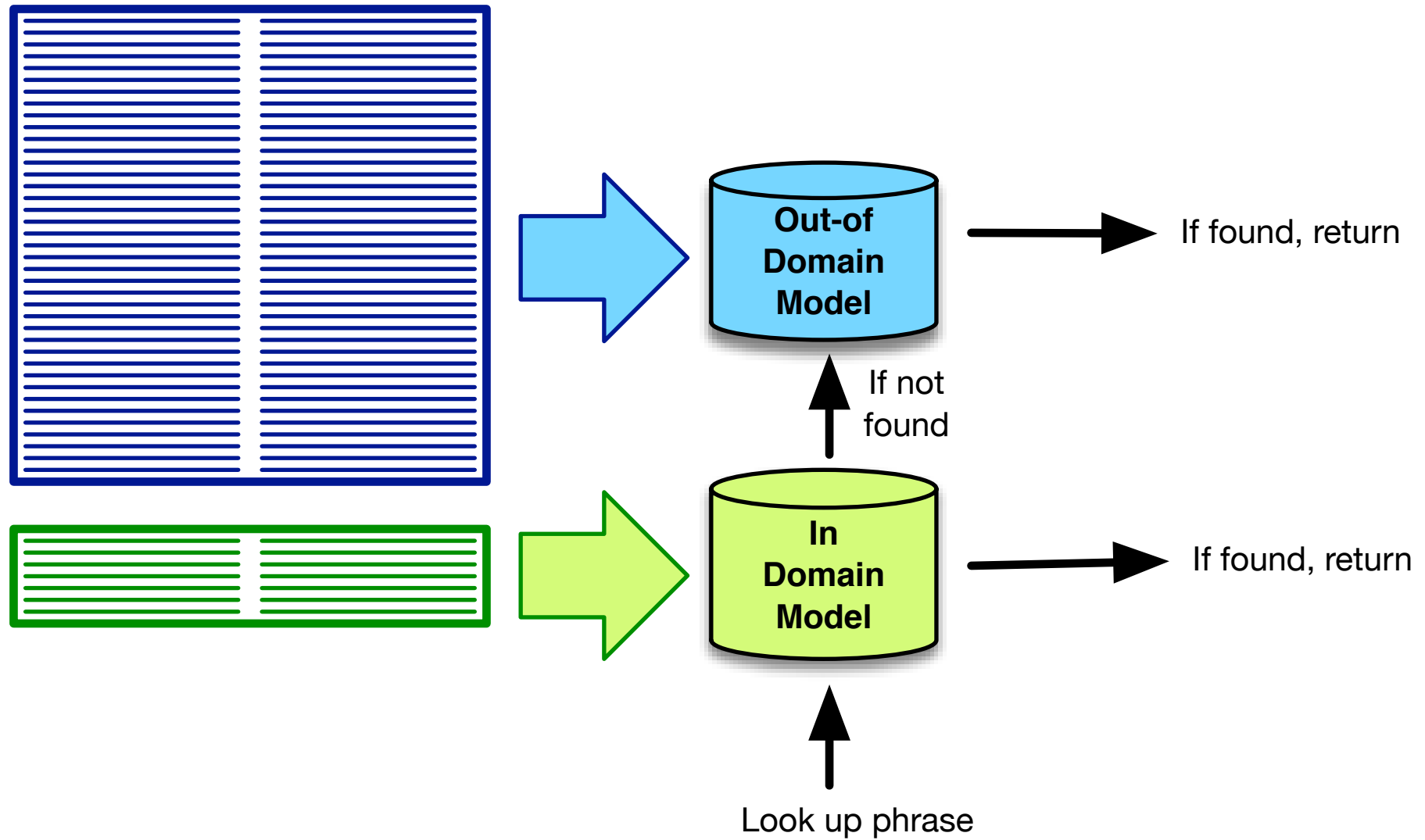
- $p_c(e|f) = \lambda_{\text{in}}p_{\text{in}}(e|f) + \lambda_{\text{out}}p_{\text{out}}(e|f)$
- Quite successful for language modelling

Multiple Models

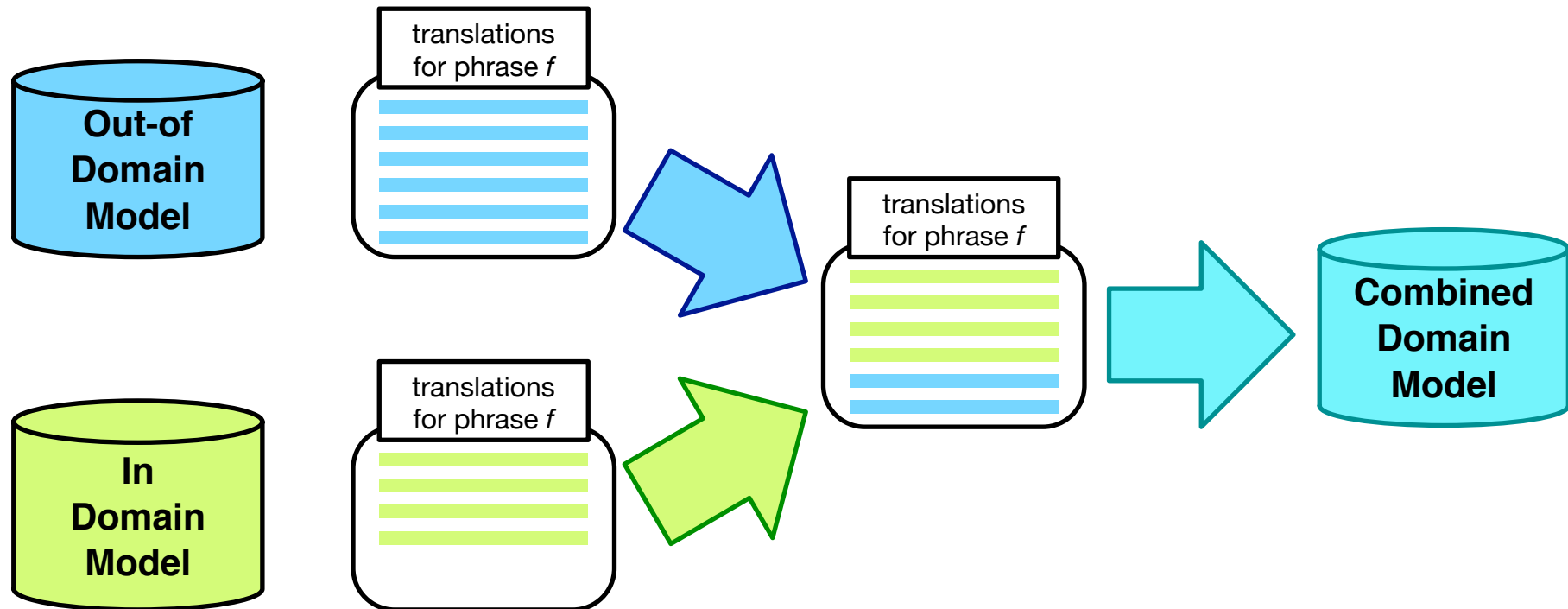


- Multiple models \rightarrow multiple feature functions

Backoff

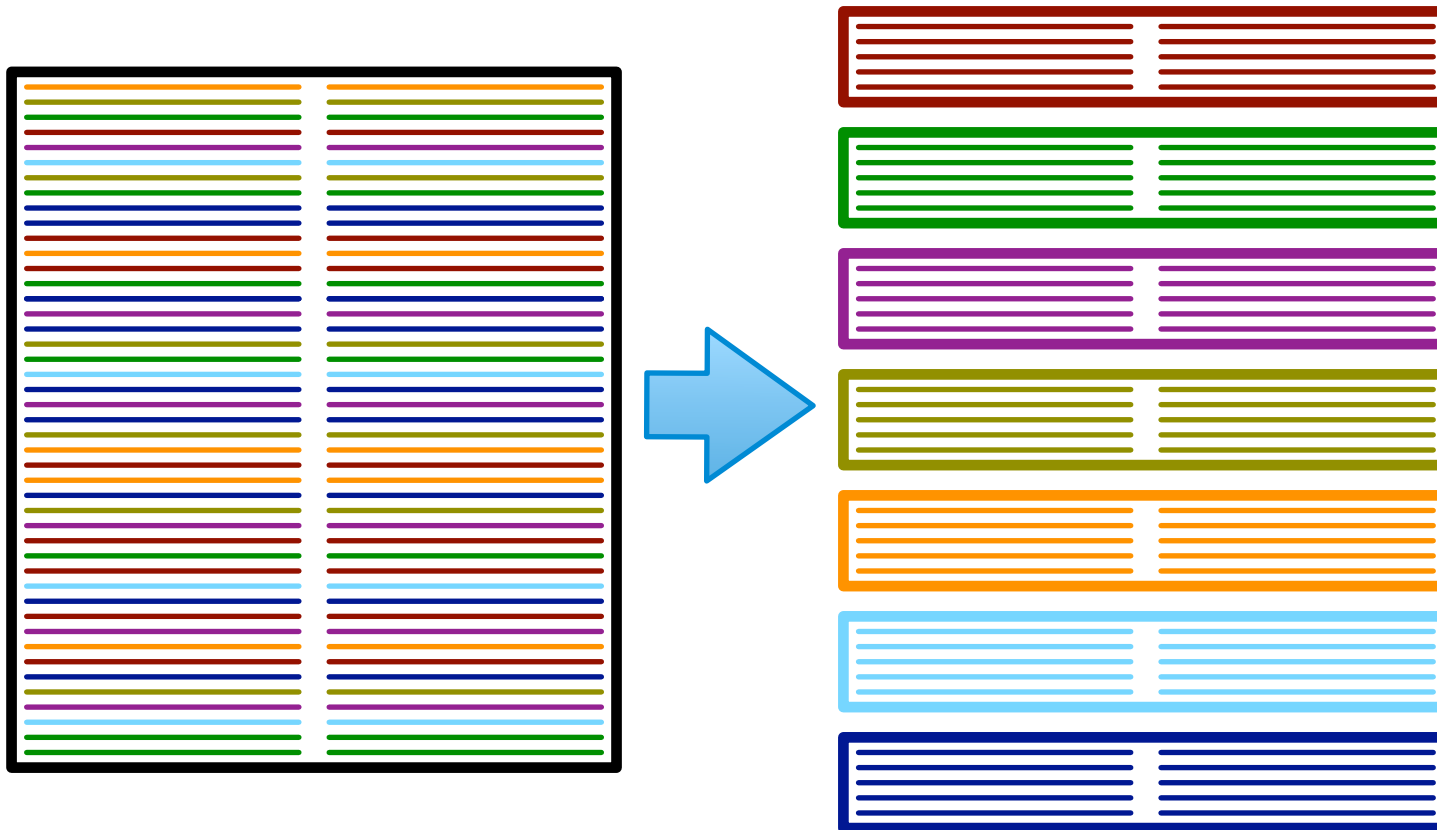


Fill-Up



- Use translation options from in-domain table
- Fill up with additional options from out-of-domain table

Topic Models

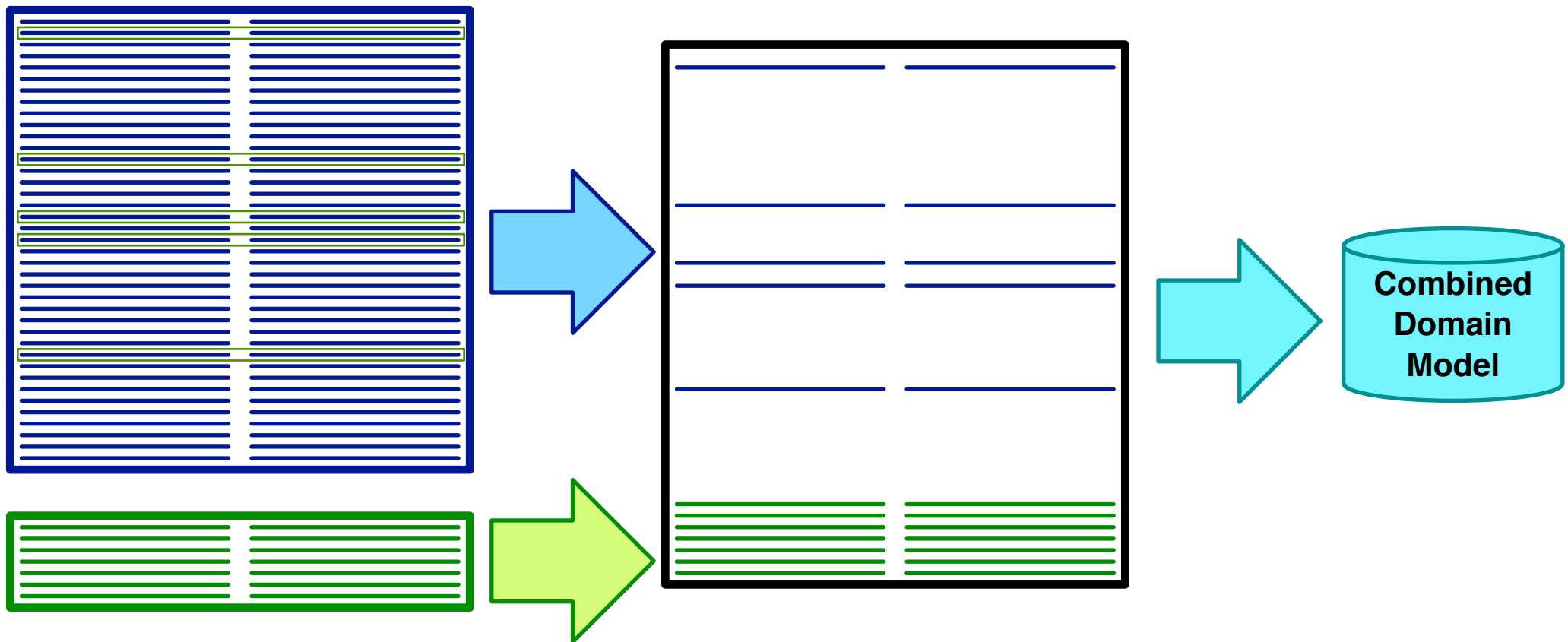


- Cluster corpus by topic — Latent Dirichlet Allocation (LDA)
- Train separate sub-models for each topic
- For input sentence, detect topic (or topic distribution)



subsampling

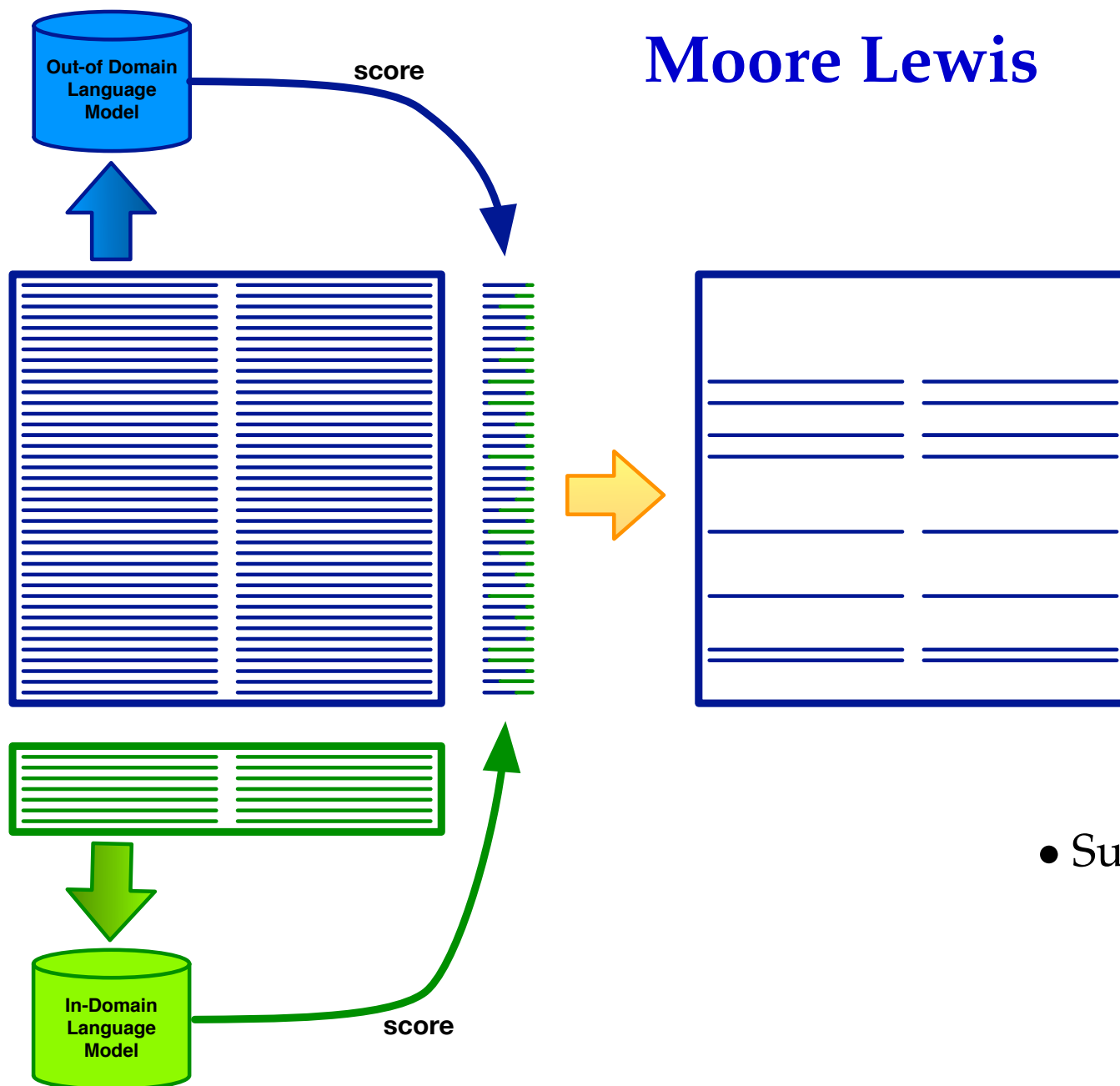
Sentence Selection



- Select out-of-domain sentence pairs that are similar to in-domain data

- Various methods
- Goal 1: Increase coverage (fill gaps)
- Goal 2: Get content with in-domain content, style, etc.

Moore Lewis



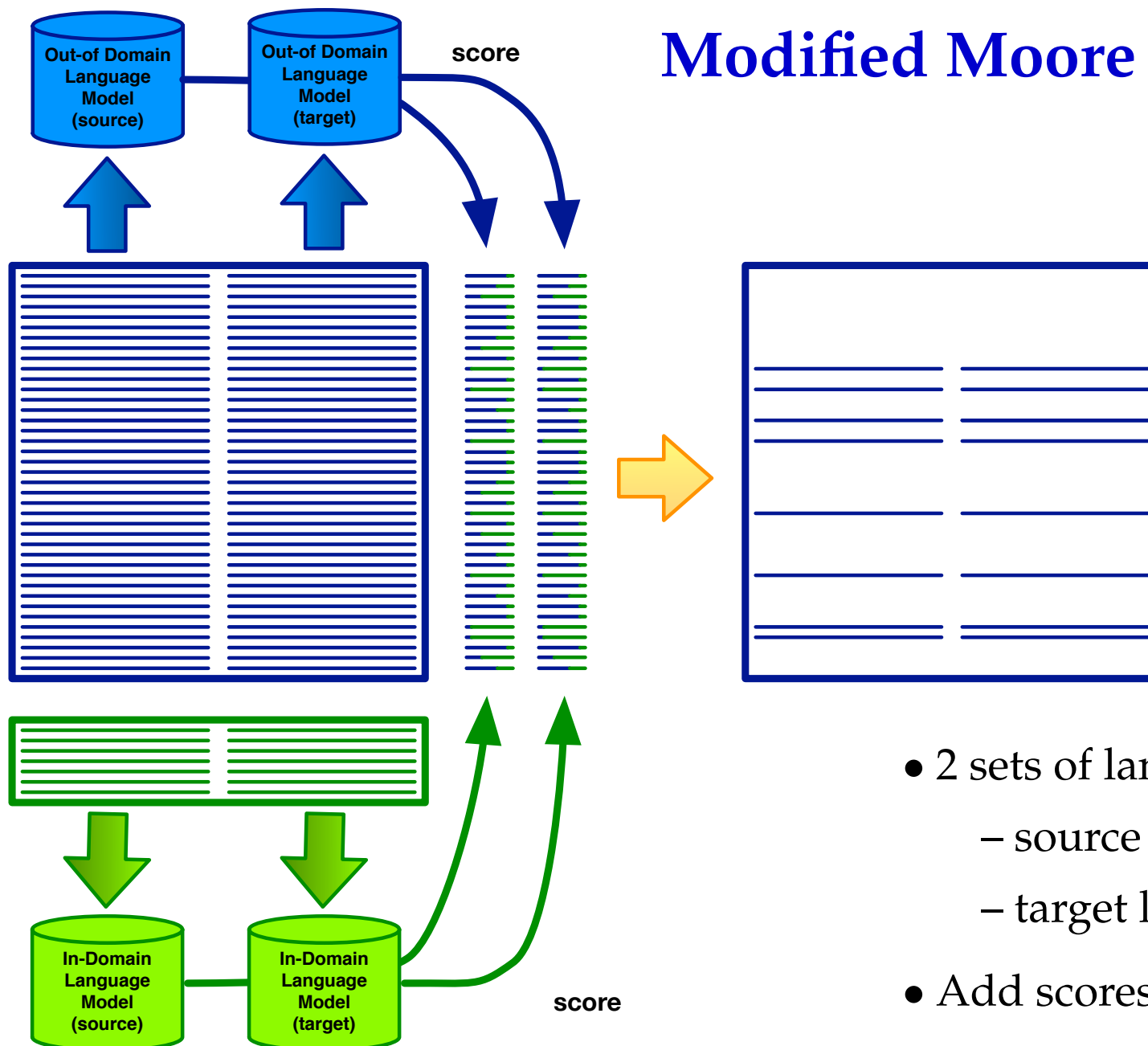
- Build language models
 - out of domain
 - in domain

- Score each sentence

- Sub-select sentence pairs with

$$p_{\text{IN}}(f) - p_{\text{OUT}}(f) > \tau$$

Modified Moore Lewis



- 2 sets of language models
 - source language
 - target language
- Add scores

Subsampling with POS

- Replace rare words with part-of-speech tags

an earthquake in Port-au-Prince



an earthquake in NNP

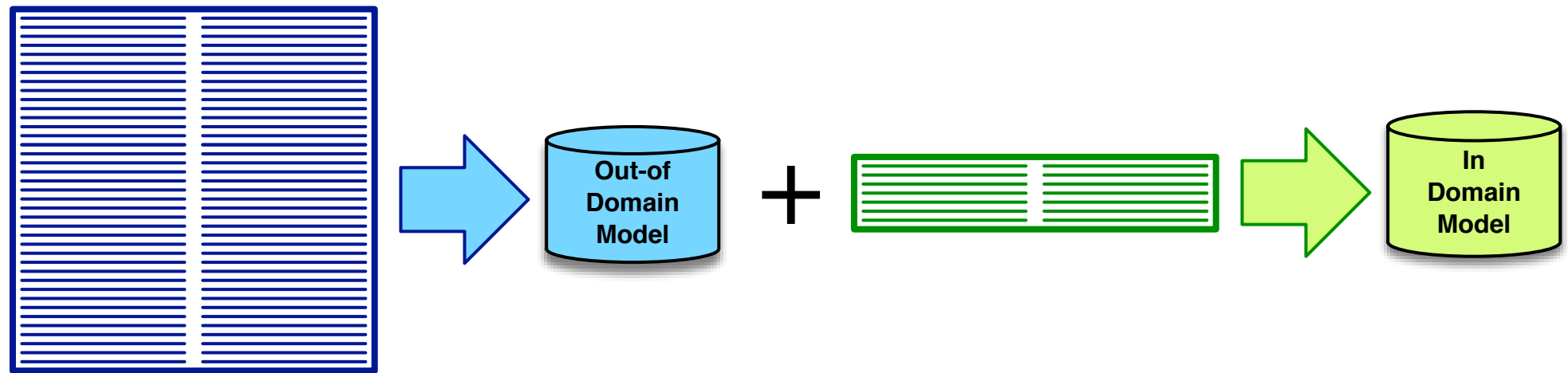
- Works better [Axelrod et al., WMT2015]
- Is it all about style, not key terminology?

Still Hard Problems

- How related are domains?
- Is corpus X useful for my system?
- What text properties matter?

neural adaptation

Fine Tuning



- First train system on out-of-domain data (or: all available data)
- Stop at convergence
- Then, continue training on in-domain data
- Successful even for fine tuning on 1 sentence pair [Farajian et al., WMT 2017]

- Given: sets of corpora with known domain
- Task: translate sentence of known domain
 - training: add domain token, say [SPORTS], to each source sentence
 - testing: add domain token to input
- Task: translate sentence of unknown domain
 - training: learn separate models for each domain
 - testing: predict domain of sentence, weight ensemble of domain-models

- Goal: Give more weight to in-domain data
- Solution: Duplicate in-domain data n times when merging
- But duplication factor not clear

Instance Weighting

- For each sentence pair, compute domain-relatedness score (0–1)
— could use something like Modified Moore-Lewis
- During training: scale learning rate based on this number

[Chen et al., NMT 2017]

questions?