



MIDEIND



A Warm Start and a Clean Crawled Corpus – A Recipe for Good Language Models

Vésteinn Snæbjarnarson, Haukur Barri Símonarson,
Pétur Orri Ragnarsson, Svanhvít Lilja Ingólfssdóttir,
Vilhjálmur Þorsteinsson & Hafsteinn Einarsson

MIDEIND EHF and UNIVERSITY OF ICELAND
SCHOOL OF ENGINEERING AND NATURAL SCIENCES



Summary and Main Contributions

Summary

- Developed a set of language models for Icelandic
 - both from scratch and based on a multilingual model
- Compared performance on downstream tasks across different base corpora and models

Contributions

- a) Language models for Icelandic
- b) Adaptations of LM for Icelandic with SOTA results
- c) Icelandic WinoGrande dataset
- d) Icelandic Common Crawl Corpus (IC3)

TLD (.is) Targeted Scraping of the Common Crawl



Icelandic is a medium resource language - how can we get more data?

The Common Crawl hosts a searchable index

- a) Search for ***.is** – response contains byte ranges in dumps (huge files)
- b) Fetch only data in byte ranges
- c) Clean up and deduplicate data

Icelandic Corpora



**News and legal
IGC
8.2 GB**

**The Icelandic
Common Crawl
Corpus
IC3
4.9 GB**

**News
2015-2021
456 MB**

**Medical
33 MB**

**Public Domain
Ebooks
14 MB**

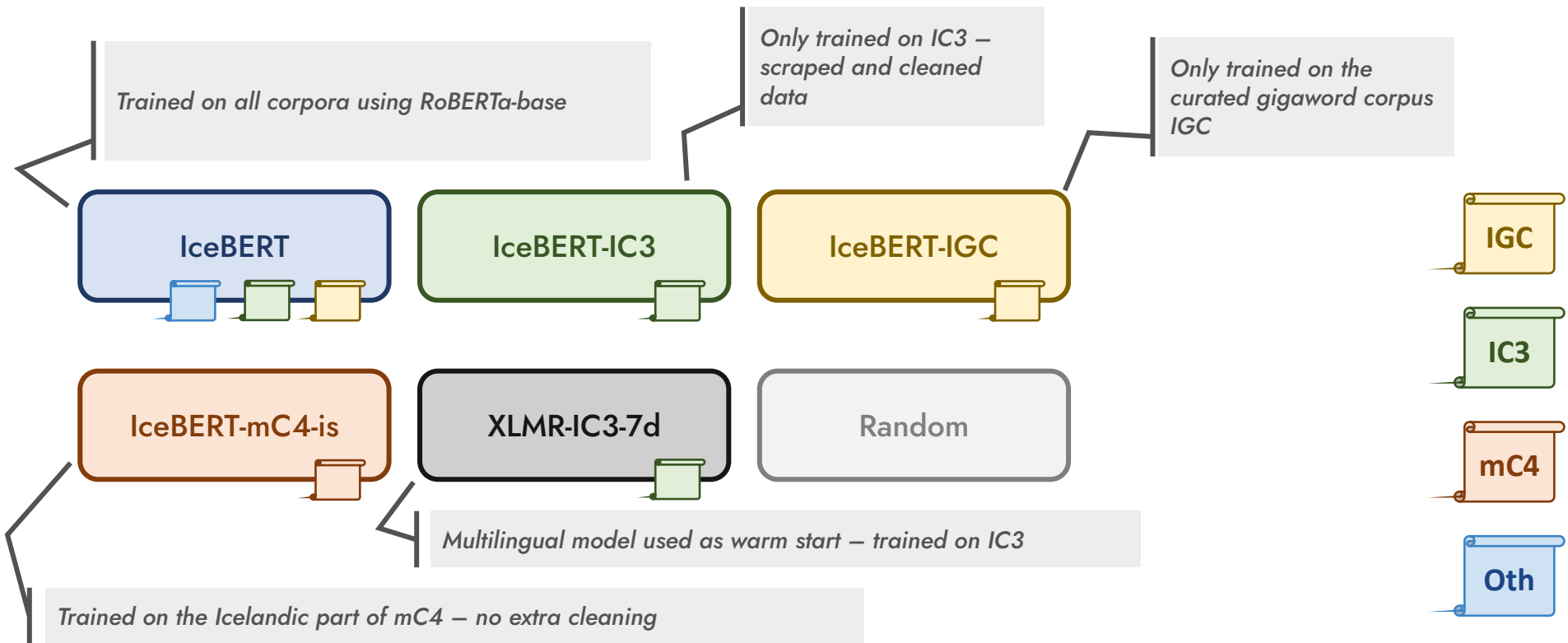
**Sagas
9 MB**

**Icelandic part of
mC4
8.0 GB**

X

**Student
thesis
2.2 GB**

Pretraining LMs for Icelandic



Part of Speech



n: noun
v: feminine
e: singular
n: nominative
g: definite article
s: proper

Part-of-speech (POS) tagging

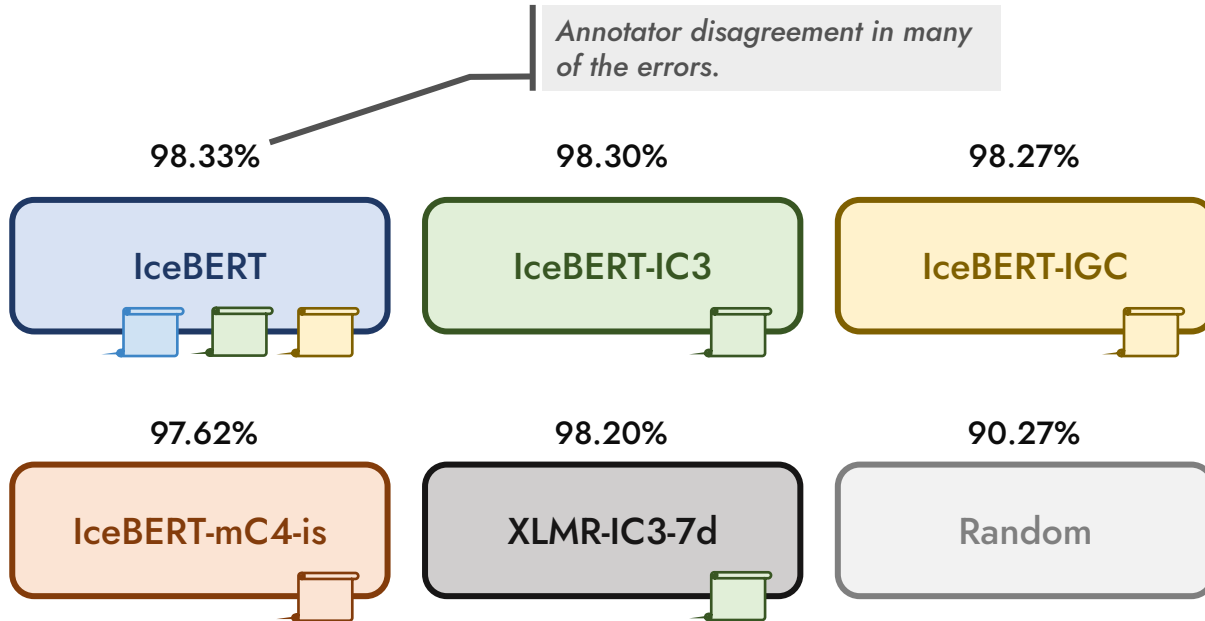
- Word type, declension, tense etc.
- Combined tags taken apart – multi label task
 - (word type, gender, case, etc)
- **Example:** *Gamla konan ók grænni rútu*

Data

- **MIM-Gold** (*Loftsson et al. 2010, Developing a PoS-tagged corpus using existing tools*)

```
<s n="2">  
  <w lemma="orð" type="nvengs">Orðin</w>  
  <w lemma="vera" type="sfg3fb">voru</w>  
  <w lemma="skrifa" type="sbgven">skrifuð</w>  
  <w lemma="með" type="aþ">með</w>  
  <w lemma="blýantur" type="nkeþ">blýanti</w>  
  <w lemma="á" type="ao">á</w>  
  <w lemma="blað" type="nheo">blað</w>  
  <w lemma="sem" type="ct">sem</w>  
  <w lemma="leggja" type="spghen">lagt</w>  
  <w lemma="vera" type="sfg3eþ">var</w>  
  <w lemma="ofan" type="aa">ofan</w>  
  <w lemma="á" type="ao">á</w>  
  <w lemma="lík" type="nheog">líkið</w>  
  <c type="punctuation">.</c>  
  <w lemma="þrír" type="tfhfn">Þrjú</w>  
  <w lemma="orð" type="nhfn">orð</w>  
  <c type="punctuation">,</c>  
  <w lemma="óskiljanlegur" type="lhfn">óskiljanleg</w>  
  <w lemma="erlendur" type="nkeþ-s">Erlendi</w>  
  <c type="punctuation">.</c>  
</s>
```

Part of Speech - accuracy



Named Entity Recognition

Groups

PERSON
LOCATION
ORGANIZATION
MISCELLANEOUS
DATE
TIME
MONEY
PERCENT



Named entity recognition (NER)

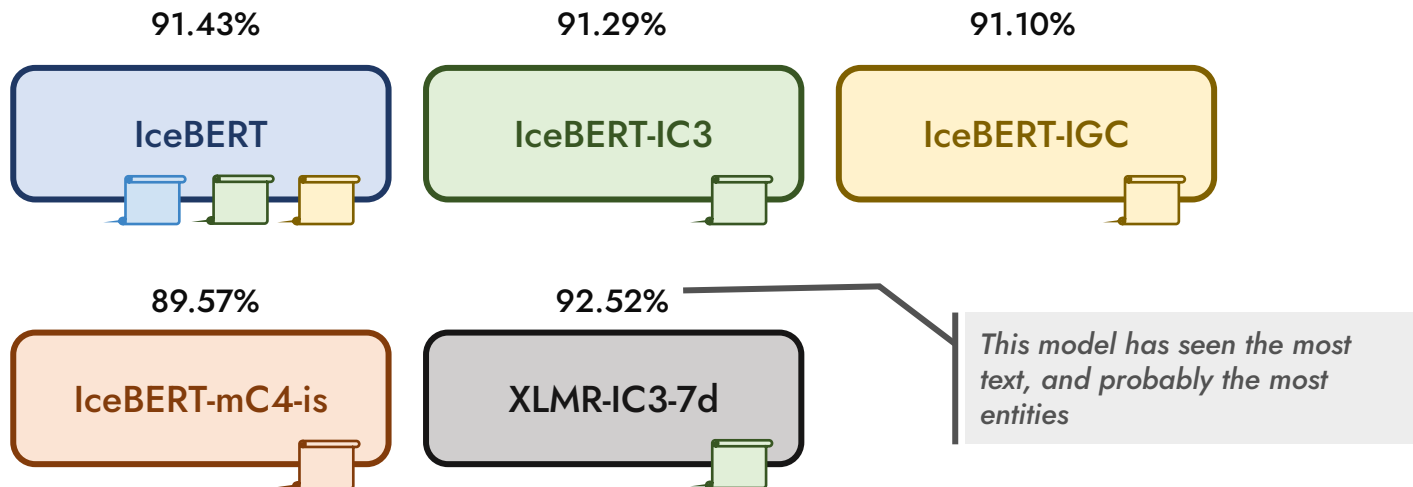
- The task of labelling entities in text
- **Example:** *Katrín Jakobsdóttir* (PER) vinnur í *Stjórnarráðinu* (ORG)

Data

- **MIM-GOLD-NER**
- Svanhvít Lilja Ingólfssdóttir, Ásmundur Alma Guðjónsson, Hrafn Loftsson

```
Vel 0
gerð 0
og 0
rómantísk 0
með 0
þeim 0
Zach B-Person
Braff I-Person
( 0
« 0
Scrubs B-Miscellaneous
» 0
, 0
« 0
Garden B-Miscellaneous
State I-Miscellaneous
» 0
) 0
, 0
Rachel B-Person
Bilson I-Person
( 0
« 0
O.C. B-Miscellaneous
» 0
þættirnir 0
) 0
ofl. 0
```


NER – F1-score



IGC

IC3

mC4

Oth

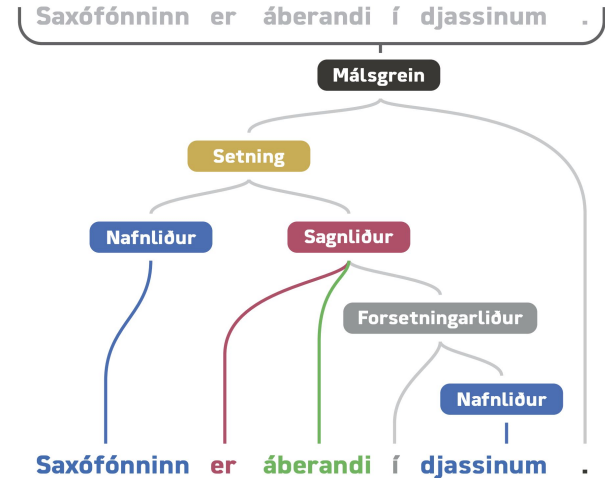


Constituency parsing

- Labeling sentence constituents

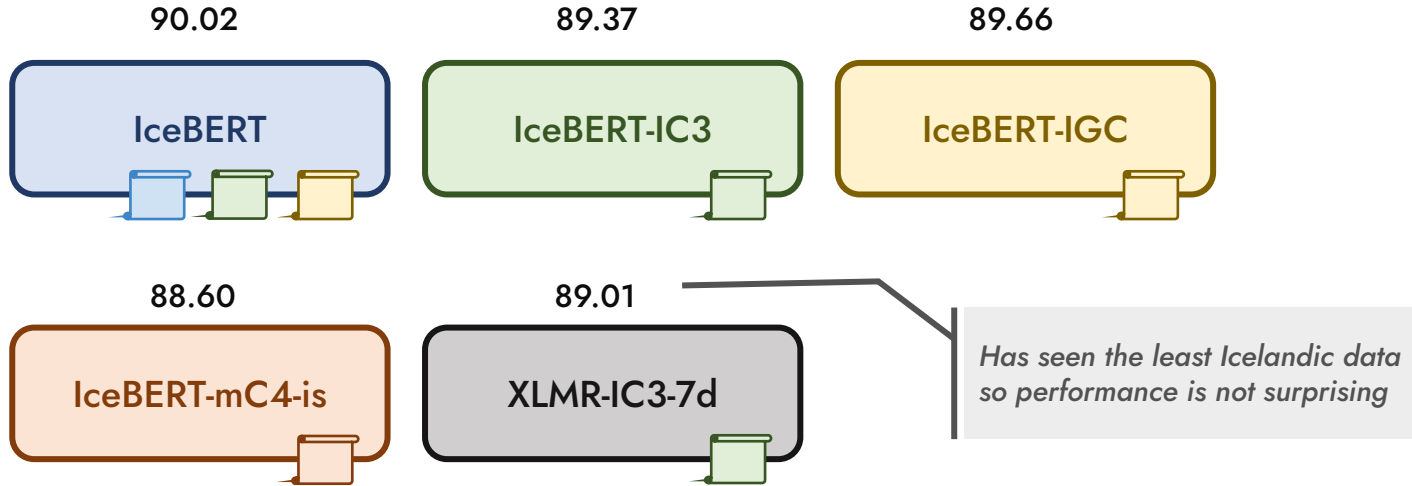
Data

- GreynirCorpus (Miðeind)



```
(S0 (S-PREFIX (C (st En (lemma en))))  
  (S-QUE (NP-OBJ (fn_et_pf_hk hvað (lemma hver)))  
    (IP (VP (VP (so_1_pf_fh_p3_et_nt_gm gerir (lemma gera)))  
      (NP-SUBJ (no_et_nf_kv_kr ríkisstjórnin (lemma ríkisstjórn)))  
      (ADVP-DATE-REL (ao pá (lemma pá))))))  
    (grm ?))))
```

Parsing – F1-score



Grammatical Error Detection



Grammatical error detection

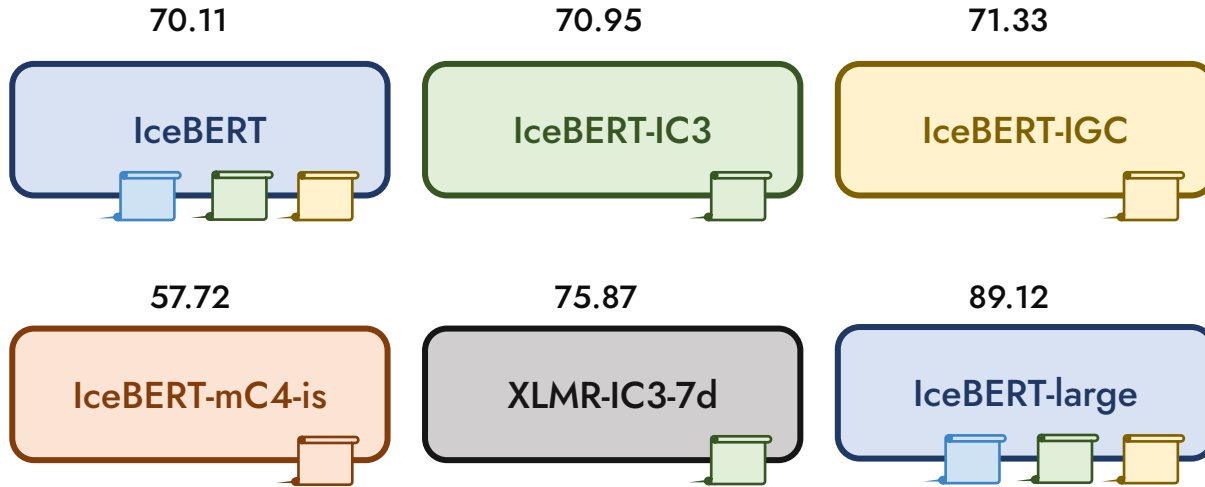
- Spelling
- Grammar
- Style

Data

- **Icelandic Error Corpus (IceEC)**
- Anton Karl Ingason, Lilja Björk Stefánsdóttir, Þórunn Arnardóttir, Xindan Xu

```
<s n="18">
  <w>Vissulega</w>
  <w>hafa</w>
  <w>kosningaloforð</w>
  <w>verið</w>
  <w>gefin</w>
  <w>um</w>
  <w>að</w>
  <w>auka</w>
  <revision id="10">
    <original><w>íslensku</w><w>kennslu</w></original>
    <corrected><w>íslenskukennslu</w></corrected>
    <errors>
      <error xtype="split-compound" eid="0"/>
    </errors>
  </revision>
  <w>fyrir</w>
  <w>innflytjendur</w>
<c>.</c>
</s>
```

Grammatical Error Detection - F1-score



Largest model is by far the best for this task





Coreference resolution

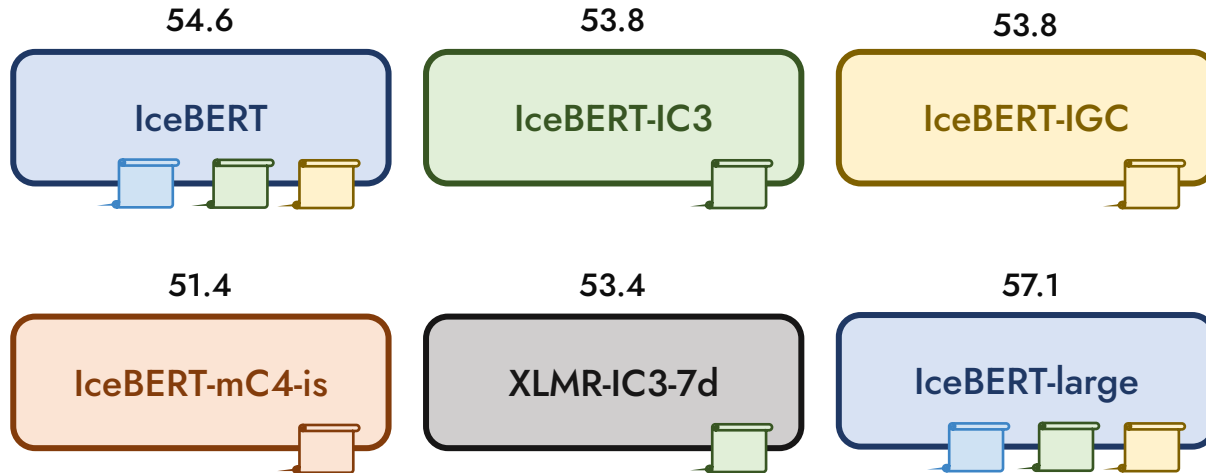
- **Example:** The **cup** could not fit into the **suitcase** because **it** was too big.

Data

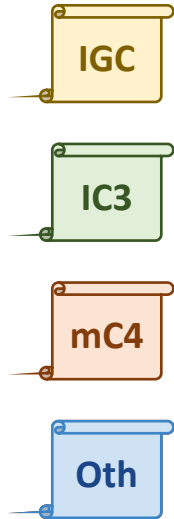
- We translated a part of the WinoGrande corpus and adapted it for Icelandic

```
{"qID": "3JVP4ZJHDPQH00KHNK0S0B899NH0I8-2", "sentence": "Þorlákína sagði Agnetu að hún kæmist ekki út að borða vegna ristilbólgu. _ sýndi því skilning.", "option1": "Þorlákína", "option2": "Agneta", "answer": "2"}
{"qID": "3JVP4ZJHDPQH00KHNK0S0B899NH0I8-2", "sentence": "Þorlákína sagði Agnetu að hún kæmist ekki út að borða vegna ristilbólgu. _ var afsökuð.", "option1": "Agneta", "option2": "Þorlákína", "answer": "2"}
{"qID": "3M47JKRKCXZJD5UJGN4IKNOMEAQ86S-2", "sentence": "Jón elskaði að borða sushi en Bjarti fannst það ógeðslegt. _ pantaði salat í matinn.", "option1": "Jón", "option2": "Bjartur", "answer": "2"}
{"qID": "32TMVRKDGPKCS7PCHXN3FN5GBLZ842-1", "sentence": "Hurðin var háværi en glugginn því hjarir _ voru betur smurðar.", "option1": "hurðarinnar", "option2": "gluggans", "answer": "2"}
{"qID": "32TMVRKDGPKCS7PCHXN3FN5GBLZ842-2", "sentence": "Hurðin var háværi en glugginn því hjarir _ voru verr smurðar.", "option1": "hurðarinnar", "option2": "gluggans", "answer": "1"}
{"qID": "306996CF6YYM26Q0XROB8RBF8LL1BN-1", "sentence": "Hann setti skýringu við myndina sína og reyndi að nota flókið mál en þurfti að endurskrifa hana því _ var of stutt.", "option1": "skýringin", "option2": "málið", "answer": "1"}
{"qID": "3PCPFX4U400L22NQ3AOM0KMZ1K3FQB-2", "sentence": "Vaka þurfti aðstoð hjá Emilíönu til að hefja brjóstgjöf, en _ gat ekki aðstoðað hana.", "option1": "Vaka", "option2": "Emilíana", "answer": "1"}
```

Coreference resolution – accuracy



Even the large model does not do much better than chance



Key Results



Feasible to extract good training data for language models from the Common Crawl

Cleaning of the data is important

Similar performance in many tasks using models trained on curated vs. scraped corpora

Using a multilingual model as a warm start leads to comparable performance to training from scratch, at a fraction of the cost

The Icelandic WinoGrande dataset is a tough task and should serve well to measure future progress