

# Generating Erroneous Text for Natural Language Corrections

Teodor-Mihai Cotet

**Thesis advisor:**

Prof. dr. ing. Mihai Dascălu

# Problem

- Scarce annotated resources for the task of grammatical error correction for most of the languages
- Neural networks requires lots of samples to train -> problem even for English

# Problem

New Message

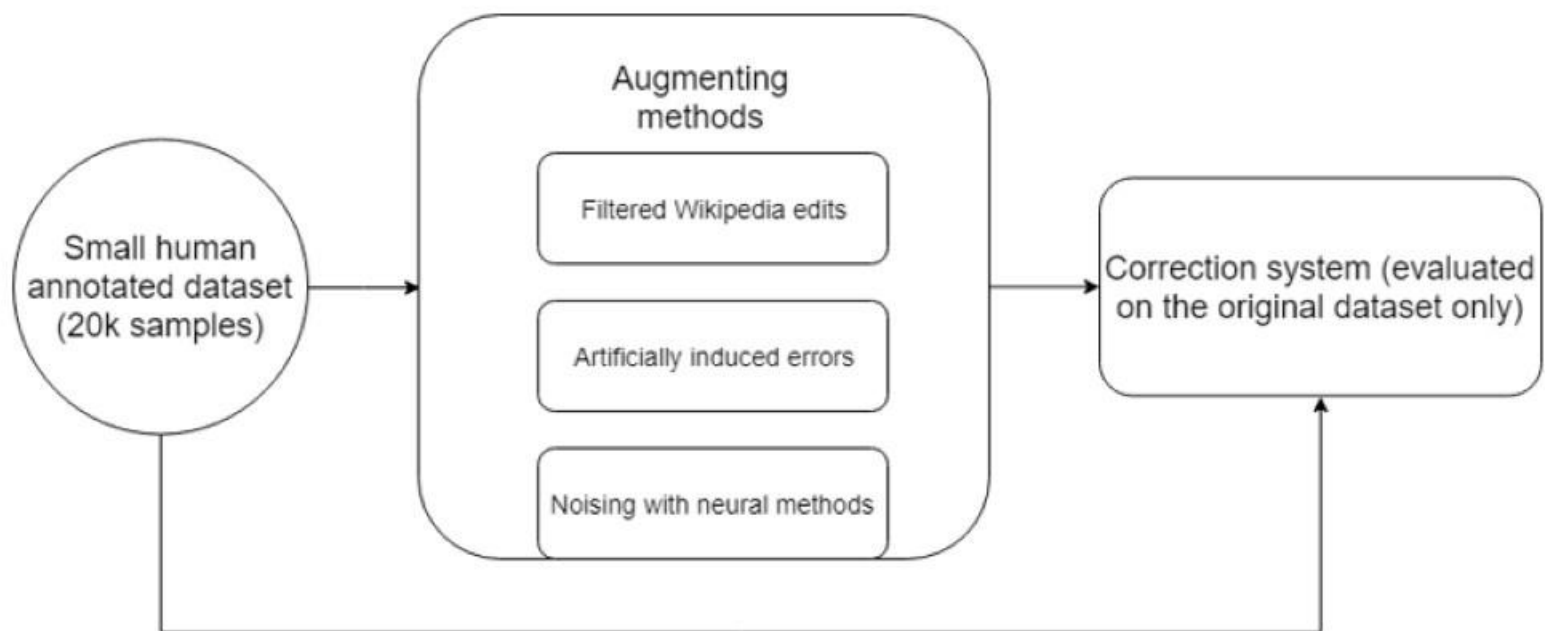


Recipients

Subject

Miniștri vroiau să îi prindă defapt cei care dintro dată gresisera pe fața.

# General pipeline



# Wikipedia edits

Before editing	After editing
<i>Este considerat</i> a fi cel mai mare dramaturg român și unul dintre cei mai importanți scriitori români.	George Călinescu <i>îl considera</i> a fi cel mai mare dramaturg român și unul dintre cei mai importanți scriitori români.
Înainte de a-și publica pamfletul în broșură, Caragiale <i>a trimis</i> [...]	Înainte de a-și publica pamfletul în broșură, Caragiale <i>trimisese</i> primul capitol [...]
[...] iar în 1892 și-a exprimat intenția de a se <i>expatria la</i> [...]	[...] și-a exprimat intenția de a se <i>exila la</i> [...]
A publicat în revista literară <i>bimensuală</i> Convorbiri [...]	A publicat în revista literară <i>bilunară</i> Convorbiri [...]
A fost numit, <i>cu decret regal</i> [...]	A fost numit, <i>prin decret regal</i> [...]
[...] <i>a avut loc prima reprezentație</i> a piesei [...]	[...] <i>premiera</i> piesei [...]
În <i>aceste</i> împrejurări și-a manifestat calitățile sale de [...]	În <i>acele</i> împrejurări și-a manifestat calitățile de [...]
[...] <i>în presa vremii</i> [...]	[...] <i>in presa contemporană</i> [...]

# Artificially induced errors

- Developed a tool to manipulate POS tags features
- Examples:
- Primul (word form):
  - RoCaseEnum: RoCaseEnum.ACC, RoCaseEnum.NOM (form of the noun)
  - RoDefiniteEnum: RoDefiniteEnum.DEF (definite form)
  - RoGenderEnum: RoGenderEnum.MASC (masculine form)
  - RoNumberEnum: RoNumberEnum.SING (number in singular form)
  - RoNumTypeEnum: RoNumTypeEnum.ORD (numeral type showing order, not comparison)
  - RoNumFormEnum: RoNumFormEnum.WORD (numeral form described through a word, not an Arab or Roman digit)



# Artificially induced errors

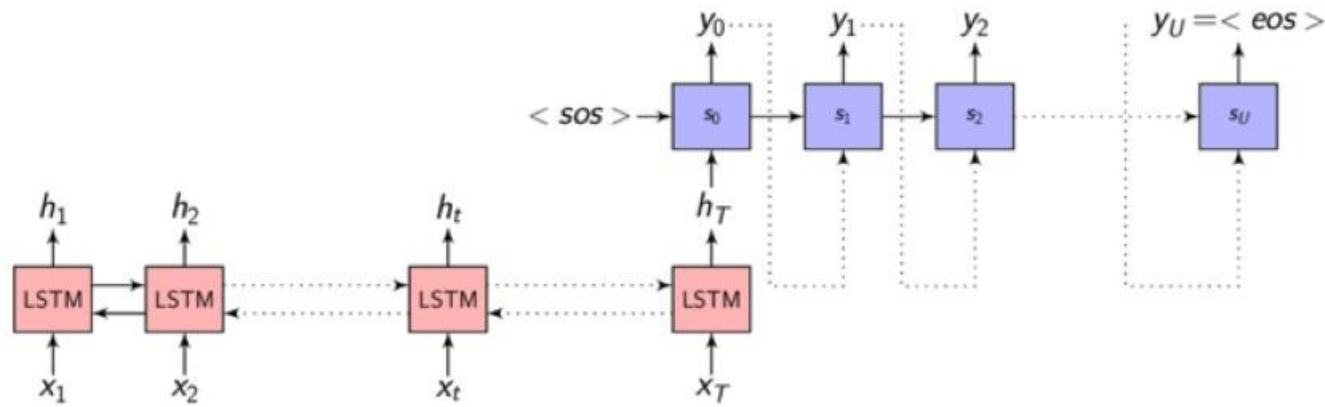
- văzuse:
  - RoMoodEnum: RoMoodEnum.IND
  - RoNumberEnum: RoNumberEnum.SING (singular)
  - RoPersonEnum: RoPersonEnum.THIRD (third person)
  - RoTenseEnum: RoTenseEnum.PQP (past perfect tense equivalent in Romanian)
  - RoVerbFormEnum: RoVerbFormEnum.FIN (general verb form)

# Artificially induced errors

Mistake type	Description	Example
<b>Tp</b>	Typos qwerty	Am mers la [mahazîn -> magazine].
<b>Diac</b>	Missing/extra diacritics	[Scoala -> Școala] din zare.
<b>Vt</b>	Verb tense	[Merg -> Am mers] la magazin ieri].
<b>Vf</b>	Verb form	Masinile sunt [prevazut -> prevazute] cu suspensii.
<b>SVA</b>	Subject-verb agreement	El [mergem -> merge] la magazin.
<b>Pre</b>	Prepositions	Acest țel oarecum a început să piară [ca și -> ca] importanță. Spre unul [din -> dintre] directorii firmei.
<b>Vir</b>	Missing/extra commas	[Astfel -> Astfel,] s-a demonstrat teorema. Soră-mea cea mare [, ->] m-a ținut odată în casă
<b>FC</b>	Spelling	[Copii -> Copiii] merg la scoala prea mult. De mult [vrolam -> voiam]. 1000 de persoane care ar putea [știi -> ști] ceva.
<b>Cr</b>	Extra/missing dash	[Dintro -> Dintr-o] greseala am plecat.
<b>Sp</b>	Spacing	deasemenea vs de asemenea, numai vs numai
<b>Gen</b>	Gender	Fata [frumos -> frumoasă] a plecat.
<b>Case</b>	Noun case	Am iesit la [pensii -> pensie]
<b>PI</b>	Strongness form of pronoun	Ea [însumi -> însăși] a câștigat.
<b>AcP</b>	Owner agreement	Băiatul [al -> a] cărui carte
<b>Cap</b>	Capitalization errors	[Ioana -> Ioana] merge la magazin.



# Inducing noise with neural methods



Modified beam search (random noising - adding  $r * \beta_{random}$  to beam candidates probabilities)

# RESULTS FOR DETECTING AND CORRECTING WORD LEVEL ERRORS

Corpora statistics

corpus	sentences/samples	Dictionary (# words)	size (MB)
inflected	650k	150k	150MB
Typos	550k	150k	152MB

# RESULTS FOR DETECTING AND CORRECTING WORD LEVEL ERRORS

Results for correction at the word level.

Task	Model	Epochs	Val Accuracy
inflected	word	33	0.519
inflected	char-GRU	46	0.518
inflected	word + sent-GRU + char-GRU	15	<b>0.610</b>
typos	word	17	0.622
typos	char-GRU	35	0.733
typos	word + sent-GRU + char-GRU	18	<b>0.738</b>

# RESULTS FOR DETECTING AND CORRECTING WORD LEVEL ERRORS

Detection statistics

Task	Model	Epochs	Val	Precision	Recall	$F_{0.5}$
Accuracy						
inflected	word	25	0.820	0.782	0.843	0.794
inflected	word + sent-GRU + char-GRU	15	<b>0.927</b>	<b>0.909</b>	0.943	0.916

## RESULTS FOR DETECTING AND CORRECTING WORD LEVEL ERRORS

Classifying by different threshold

Threshold	Precision	Recall	$F_{0.5}$
0.50	0.909	<b>0.943</b>	0.916
0.40	0.926	0.933	0.928
0.25	0.947	0.915	0.940
<b>0.10</b>	<b>0.971</b>	0.879	<b>0.951</b>

## Examples from the detection system (for threshold 0.5).

Text in bold marks incorrect initial text

Sentence	Type
— vrea/vrem oare old shatterhand să mă abată de la datorie ?	tp
vreau și eu pedepsirea/ <b>pedepsire</b> ucigașilor .	tp
datinile îmi poruncesc să rămân lângă morții mei/ <b>meu</b> pînă vor fi înmormîntați .	tp
— și cînd va/ <b>vruseră</b> fi înmormîntarea ?	tp
pentru că nu pot/ <b>putea</b> să - l urmăresc personal , sarcina va reveni altcuiva .	tp
Acum/ <b>acu'</b> , cînd era vorba de chestiuni concrete , winnetou deveni calm ca de obicei.	tp
— acesta e numele lui/el adevărat ? el lui	tp
— aveți/ <b>avuseși</b> cunoscuți prin apropiere , poate în vreun fort ?	tp
— voiam/ <b>voi</b> să ... să ne luăm după ...	tp
- am/ <b>avea</b> ghicit în să gîndul și i - am întrebat :	tp
sau aveati/ <b>avea</b> de gînd să - i atacați la întoarcere și ...	tp
trecuse mult de amiază și sam mă întîmpină intrigat/ <b>intrigată</b>	tp
se lăsă/ <b>lăsa</b> o tăcere adîncă .	tp
apoi v - ați ascuns în <b>pădurice</b> și ne - ați observat de susul , din copaci . pădurice	fp
sper că voi răspunde cu cinsteo <b>încredeții</b> pe care fratele meu roșu mi - o acordă .	fp
gîndul <b>acesta</b> îmi grăbi pasul .	fp
n - avea izbutit din primul moment , pămîntul <b>fiind</b> aici încă prea dur .	fp
celi trei cai se mai afla în pădurice .	fp
trecuse mult de amiază și sam mă întîmpină intrigată :	fp
avea , nu zic ba .	fp
zarvă nu duce la nimic	fp
gîndul acesta îmi grăbi pasul/ <b>pasii</b> .	fn
în loc de răspuns , i - am chemat pe apasi/ <b>apaș</b> la mine .	fn
așadar , ce aduce tovărășia mea : viață sau moarte/ <b>morții</b>	fn
Draga/ <b>dragă</b> de ea , frumoasa și buna fată indiană	fn
— mai tacă/ <b>tăcu</b> - ți gura cu greenhorn - ul dumitale și lasă glumele !	fn
santer poate să treacă fie peste munți/ <b>munte</b> , fie printre munți , cum îi convine .	fn



# Future work

- The focus will be on generating a good copora for Romanian using as few language specific features as possible
- Evaluate how good the corpus is and its methods of generation by implementing a correction system

Thank you!