

# Glyph embedding for NMT

B-Town and K-Dog

<https://github.com/lepidodendron/lepidodendron>

# Overview

1. Motivation
2. Implementation
3. Analysis
4. Future maneuvers

Wir begrüßen diese Entwicklung sehr.  
We very much welcome this development.  
We welcome this development. Wir, weooooo. u. occurs theooest us

Die Globalisierung hat unsere Industrien kaputt gemacht.  
Globalism has destroyed our industries.  
Globalisation has made our industries. reooo oo. oooo oo. meoo oo

Das dürfte die meisten von uns nicht überraschen.  
This is no surprise to most of us.  
That noone should not be surprised by most of us.

# Motivation

# An alternative to learned embedding



- Problems with learned embedding
  - Vocabulary (over words, ngrams, or characters) is fixed and difficult to extend
  - Update frequencies for entries extremely unbalanced (Zipfian)
  - Can never model a language fully, even on the character level ( $2^{16}$  code points)
- Using glyphs as character embedding
  - Fixed data shape instead of fixed vocabulary size
  - Input space becomes continuous and open
  - Easier to visualize some choices made by model
  - Naturally suited for logographic languages
    - 火 *fire* 炎 *fiery* 灭 *snuff* 灰 *ash* 炭 *coal*
- Glyph embedding for NMT
  - End-to-end image translation

# Implementation

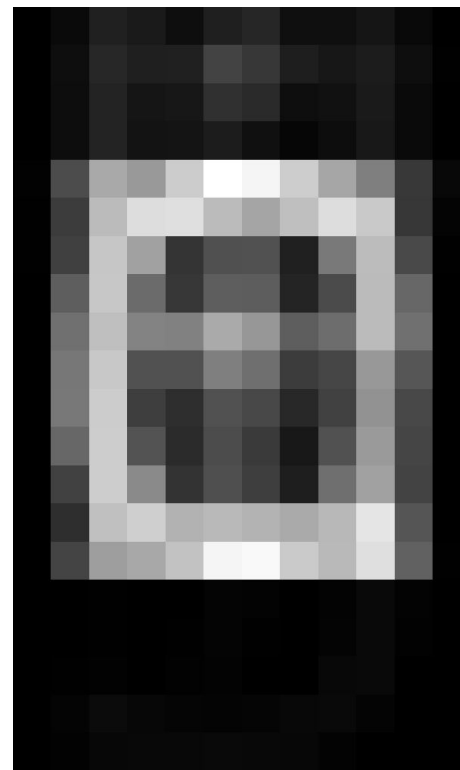
# Dataset

- [European Parliament Proceedings Parallel](#) Corpus 1996-2011
  - Paired corpus, not direct translation
- Experiments
  - Experiments conducted DE → EN
  - Took only sentences within [3,256] characters
  - 1 574 071 training instances
  - 4 096 validation instances

# Converting strings to glyphs

- [noto sans mono](#) font, size 20, rendered with [pillow](#)
- height and width fixed by the largest frequent character
  - frequent characters are the ones that cover 99.95% of the texts
- pixels in grayscale [0,255], scaled to [0,1]
-  for *unk*,  for *eos*, space for *bos*
  - *unk* is only relevant when we are not using glyphs

lang	charset	frequent chars	glyph dimension
DE	306	78	300 = 12x25
EN	293	72	240 = 12x20



Average English glyph

# Rendered inputs

source

War dies nicht der Deal, den ich bereits 2005 prognostizierte?  
Das wiederum ergibt im Spanischen responsabilidad democrática .  
1. Iran: der Fall Roxana Saberi  
Ihr Kampf um die Macht ist im Grunde ein Kampf um die Drogen.  
Dies ist von entscheidender Wichtigkeit.  
Nach der Tagesordnung folgt die Fragestunde (B5-0033/1999).  
Immerhin steht Korruption in Kamerun auf der Tagesordnung.  
Sie gründen im Kern in der Würde des Menschen.

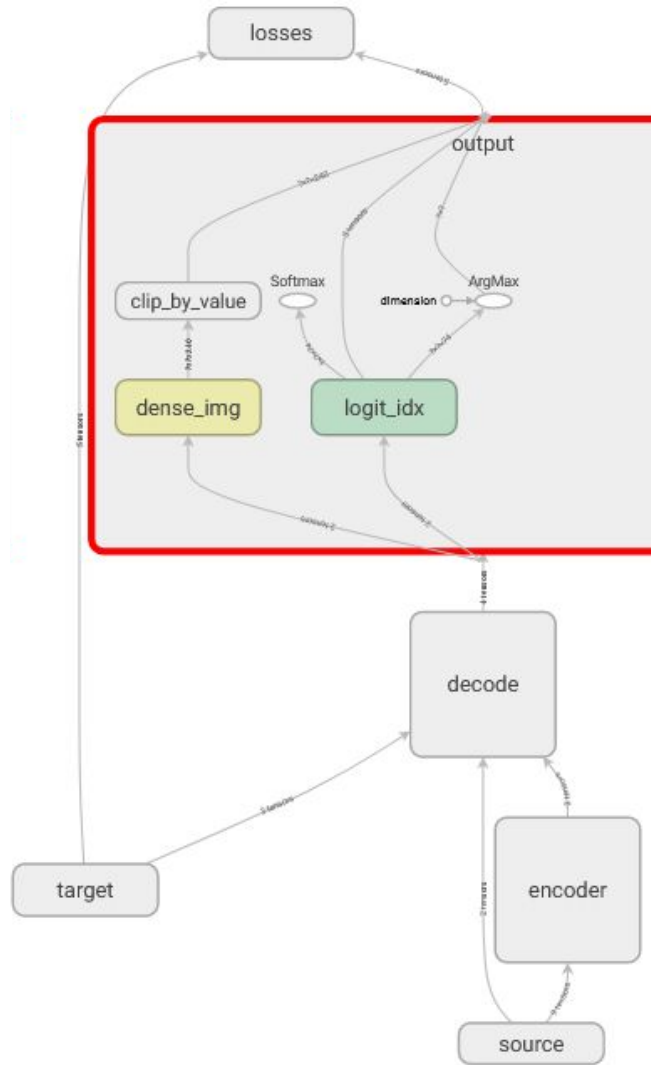
target

Indeed, was this the deal that I predicted back in 2005?  
In Spanish this gives, "responsabilidad democrática' .  
1. Iran: the case of Roxana Saberi  
Their fight for power is essentially a fight for drugs.  
This is absolutely vital.  
The next item is Question Time (B5-0033/1999).  
Corruption is a fact of life in Cameroon.  
They are, in essence, founded upon the idea of human dignity.



# A standard nmt architecture

- RNN encoder
  - 3x stacked bidirectional GRU
- RNN decoder with attention
  - 3x causal GRU, followed by
  - multi-head scaled dot-product attention, with
    - residual connection
    - layer normalization
- hidden state units: 512
- dropout: 0.1
- adam with learning rate decay
- teacher-forcing training with batch size 128





# Inputs and outputs

- **xyz** notation, **c** for character, **g** for glyph
  - **x**: encoder input
  - **y**: decoder input
  - **z**: decoder output
  - experiments: **ccc**, **cgc**, **cgg**, **ggg**
- decoder outputs
  - **z = c**: a dense layer after attention predicting chars, with softmax and XENT loss
  - **z = g**: a dense layer after attention predicting glyphs, with [0,1] clipping and MAE loss
  - whenever **z = g**, we also included a parallel char layer

# Mismatch between glyph and char predictions

Dem stimme ich zu.  
I agree with that.  
I agree with this. 



I agree.with thas.■...■.a.■.bt...■.....■ust...■.T.■ .bel■.sta.

Einige von Ihnen haben Kapitel 23 über Korruption erwähnt.  
Some of you have mentioned Chapter 23 on corruption.   
Some of you mentioned Chapter 23 on corruption. 

Some of you mentioned Chapter 23 of corruption.■ .■ as.e. ■S■...

# Mismatch between glyph and char predictions

Gibt es Einwände?  
Are there any comments?  
Are there any comments.                us.   there ar

Are there any comments?  a?be???an?ture?s??ook???  u???r there?ar

Denken Sie an Vilvoorde!  
Remember Vilvoorde!  
Think of Vilvoorde.   us.   use is          on    

Think of Vilvoorde           use.is tuturaed      an aask    

# Inference modes

mode	as autoregressive feedbacks	for bleu scores	z	y
1	g	argmax over c	g&c	g
2	g	g matched to char	g	g
3	g matched to char (and rendered)	g matched to char	g	c (g)
4	argmax over c (and rendered)	argmax over c	c	c (g)
5	probs over c (and rendered)	argmax over c	c	c (g)

- glyph-to-char matching (discretization) according to MAE
- inference terminates when the maximum step 256 is reached
- predicted strings are trimmed at the first eos

# The 5 inference modes

Aber auch dieser Punkt ist noch offen.  
But this too remains an open question.

But this issue is still open. uero. uet qq. oot qq.ms also open.

But this issue is still open. ark. ust ... ak a ms also open.

But this issue is still open. uero. uet qq. oot qq.ms also open.

But this issue is still open. uero. ue. oo. uo. oo as also open.

But this is also open to this point. uoocia. uro. ust: this is a

But this is also open to this point. uoocia. uro. ust: this is a

But this is also open. ubility. ust. . is also open. ust. . is

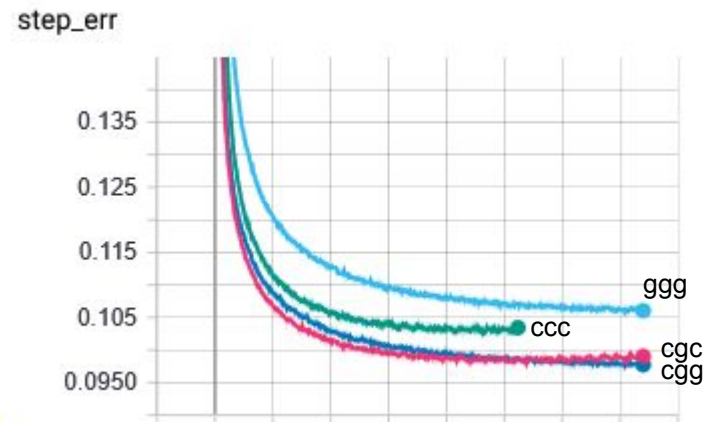
But this is also open. ubility. ust. . is also open. ust. . is

Best ieteraty, toolegene, is anount. u. u. , qq, torn again. .

Best ieteraty, toolegene is anount. . . . C. torn again. .

# Results

- `sacrebleu -tok intl`
- **ccc** overfit early



model	epochs	mode 1	mode 2	mode 3	mode 4	mode 5
ccc	30	n/a	n/a	n/a	<b>30.9</b>	20.3
cgc	60	n/a	n/a	n/a	<b>30.9</b>	22.4
cgg	60	<b>22.9</b>	<b>23.5</b>	<b>26.5</b>	<b>30.9</b>	<b>24.1</b>
ggg	60	21.9	22.2	25.3	30.2	23.0



# Analysis

# Decision points

Freizügigkeit für Personen

Free movement of persons

Freedom of movement for people's ■■■■■■■■■■ freedom of movement

Free movement of persons ■■■■■■■■■■ of movement of persons ■■■■■■■■■■

cgg: Free om of movement for people ■■■■■■■■■■ free om on migpeers

ggg: Free movement of persons ■■■■■■■■■■ y brto ovc of tucomaa o. ■■■■■■■■■■ .ortei ■■■■■■■■■■ (■)

Wir begrüßen diese Entwicklung sehr.

We very much welcome this development.

We welcome this development. ■■■■■■■■■■ We welcome this development. ■■■■■■■■■■ occurs the most us

We very much welcome this development. ■■■■■■■■■■ of these. ■■■■■■■■■■ uels. ■■■■■■■■■■

cgg: We velcome this development. ■■■■■■■■■■ SS werkine. ■■■■■■■■■■ 0 wncurs toisgest ■■■■■■■■■■ .s

ggg: We very much welcome this development. ■■■■■■■■■■ of these. ■■■■■■■■■■ eth ■■■■■■■■■■ uels. ■■■■■■■■■■

# Decision points

Das ist jedoch nicht geschehen.

This is not, however, what has happened.

That has not happened, however. **deon. uator. d. andeas. methid.**

That has not been done. of the accounts, none

**cgg:** That has not happened. however.■Nesn.■ustir.■ andenl. Tuthuld.■

**ggg:** That has not been drammed. .af the sucomda. Aask Aitomis. hot i

Das dürfte die meisten von uns nicht überraschen.

This is no surprise to most of us.

That notice should not be surprised by most of us. [REDACTED], [REDACTED]

That ~~is~~ should not be surprised by most of us. 1000000.00 000

**cgg:** That itdris should not be surprised by most of us. as bod. au t

**ggg:** That stvses should not be surprised by most of us. .er an tun

# Imagined words

Aber diese Denkungsart ist völlig falsch!

How misguided we are!

But this ~~manoul~~ is completely wrong. ~~us us an aorooo. aooo aa. a~~

But this is completely wrong. ~~uete. uaitia. usitooo the oocessin~~

**cgg:** But this manoul is completely wrong. ~~??n aurr??~~

**ggg:** But this ms completely wrong ~~ueti? uasti? usts p the fittestin~~

Aber Sie werden sie bekommen.

Rest assured that you will receive a reply in due course.

But you will ~~ooooo~~ them. ~~ooooooo. ooooo. ooooo. on. u~~

But you will ~~optone~~ them. ~~usite. ooisons. of them. of the oos~~

**cgg:** But you will brcopp them. ~~...n.... .esn. ... .asn Cner.yu.~~

**ggg:** But you will gnli e them. ~~uetle. .tunlens. if them. .tf the Uuc~~

# Synonym struggles

Das ist unangenehm.  
That is embarrassing.  
That is unclased. I. I. I. is uncecebul. a. U. hI  
This is anpleased. of ooto. oones. of the oourity. of t

cgg: That is unclased. uTE H is uncecembul. a. T

ggg: This is anpleasid. if tonts. titkes. if the Etturity. if t

Das ist ein wichtiger Hinweis.  
It is important that that should be said.  
Thas is an important ooder. uo. oo. ust oot. o. oo. usto. to. ande  
This is an important comment. between to be and ual. ieco. to o

cgg: This is an important polter. unt. a. at urt uamr usti B2 rnde

ggg: This is an important plmment. eswiudh. to te ane uali ient. utnm

# Synonym struggles

Dies ist von entscheidender Wichtigkeit.

This is absolutely vital.

This is critical. ■ustantod ■ssod ■g n. ■uet ■o. ■o is ■ssodful. ■o.

This is crucial. ■ssod. ■ssodssioo. ■o ■ootod ■ano and the ■ssodt

cgg: This is vriticall ■ stunt.. ■.. ■.. ■.. ■uet ■.e. ■ua is auveful. ■.s.

ggg: This is crucial. ■ieds. ■ wutkrios. ■ if tuttoe ■.ann. and ihe Eutult

Wir haben einen völlig unwahren Text angenommen.

The text which we have approved is totally wrong.

We have adopted a completely untrastreated text. ■ssopeoooo. ■ssodoo.

We have adopted a completely unrraoiable text. ■uais. ■uals ■ssodoo

cgg: We have adopted a tompletely untrubteated text. ■ aspea... ■.a....

ggg: We have adopted a completely untransable text. ■uats ■uslt uptu t

# Is this the end?

Die Kommission reagierte sofort.  
The Commission reacted very promptly.  
The Commission responded immediately.  
The Commission immediately responded to the

cgg: The Commission reapeded immediately. wint ne... .ezn. .

ggg: The Commission rmmediately responded to the nttenvattat. .if tha

Ich hoffe sehr, dass wir ein ...  
I hope very much that we will maintain a ...  
I very much hope that we andeal ....  
I very much hope that we oneurs Issutination. I the

cgg: I very much hope that we wnderl ..t .n u ....chrk. e

ggg: I very much hope that we .neur..f.u....li....I .lsosee..he....s

# Is this the end?

Das ist die eine Seite.

That is what is happening on the one hand.

That is one side of the `oooooooo` `uot` `u` `ure` `uot` `u` `neon` `note` `uot`

cgg: That is one side of the `ppc....` `s` `..` `urt` `urt` `riane` `srth` `ust`

Es ist auch leicht zu erklären, warum.

It is also easy to explain why.

It is also easy to explain why `toou` `iosk` `.` `in` `in` `turt` `u` `an`

cgg: It is also easy to explain why `toou` `iosk` `.` `in` `in` `turt` `u` `an`



# Hour format

(Die Sitzung wird um 16.25 Uhr geschlossen.)

(The sitting was closed at 4.25 p.m.)

(The sitting was closed at 16.25 p.m.)

(The sitting was closed at 4.25 p.m.)

cgg: (The sitting was closed at 4.25 p.m.)

ggg: (The sitting was closed at 4.25 p.m.)

Die Abstimmung findet heute um 18.30 Uhr statt.

The vote will take place today at 6.30 p.m.

The vote will take place today at 6.30 p.m.

The vote will take place today at 6.30 p.m.

cgg: The vote will take place today at 6.30 p.m.

ggg: The vote will take place today at 6.30 p.m.

# Unknown chars get correct glyphs

Bericht Theato (A5-0090/1999)

Theato report (A5-0090/1999)

Theato report (A5-0090/1999)

cgg: Theato report (A5-00901999) S E S )

Anfragen an Herrn Byrne

Questions to Commissioner Byrne

Questions to Mr Byrne

cgg: Questions to Cr Byrne

Simbabwe war einst ein blühendes Land und könnte es wieder sein.

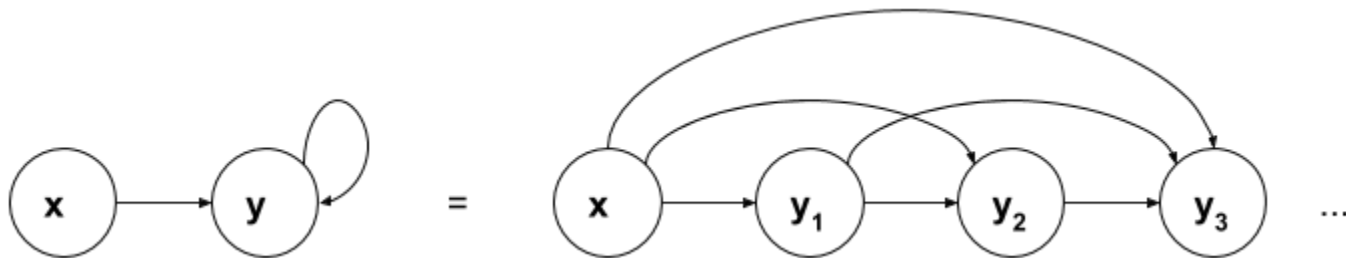
Zimbabwe was once a flourishing country and it could be again.

Zimbabwe was once a prouoitiing country and could be recorved a

cgg: Zimbabwe was once a pron ri ging country and could be repterved.a



# Excursion: autoregressive seq2seq models



- $p(y | x) = p(y_1 | x) p(y_2 | y_1, x) p(y_3 | y_2, y_1, x) \dots$ 
  - Each  $y_n$  is a discrete random variable over the target vocabulary  $T$
  - The whole search space is  $T^n$
- But glyphs are not predicted as a random variable
  - A glyph-predicting decoder is not truly autoregressive
  - Unless we perform discretization (mode 3)

# A second look at the results

- Stochastic modeling works better (modes 3&4), however
  - fuzzy glyphs are less problematic than fuzzy chars (mode 5)
- The decoder benefits from using glyphs
  - a non-autoregressive model with convolution may benefit more

model	mode 1	mode 2	mode 3	mode 4	mode 5
<b>ccc</b>	n/a	n/a	n/a	30.9	20.3
<b>cgc</b>	n/a	n/a	n/a	30.9	22.4
<b>cgg</b>	22.9	23.5	26.5	30.9	24.1
<b>ggg</b>	21.9	22.2	25.3	30.2	23.0

# Future maneuvers

采薇采薇薇亦作止  
曰歸曰歸歲亦莫止  
靡室靡家玁狁之故  
不遑啟居玁狁之故  
采薇采薇薇亦柔止  
曰歸曰歸心亦憂止  
憂心烈烈載飢載渴  
我戍未定靡使歸聘  
采薇采薇薇亦剛止  
曰歸曰歸歲亦陽止  
王事靡盬不遑啟處  
憂心孔疚我行不來  
彼爾維何維常之華  
彼路斯何君子之車  
戎車既駕四牡業業  
豈敢定居一月三捷  
駕彼四牡四牡騤騤  
君子所依小人所腓  
四牡翼翼象弭魚服  
豈不日戒玁狁孔棘  
昔我往矣楊柳依依  
今我來思雨雪霏霏  
行道遲遲載渴載飢  
我心傷悲莫知我哀

# Logographic language translation

- Chinese to English translation with [UNCorpus](#)
- Experiments
  - We take only sentences within [3,128] characters
  - 7 392 227 training instances
  - 4 096 validation instances
  - To compare **cgg** and **ggg**

lang	charset	frequent chars	glyph dimension
ZH	6357	2681	500 = 20x25
EN	714	80	240 = 12x20



# Directions for future works

- CNN
  - convolution over rendered image
  - transposed convolution for non-autoregressive prediction
- Without character boundaries
  - simply an image instead of a sequence of glyphs
  - postnet for word (piece) prediction with CTC loss
- Multiple fonts and typefaces
  - combinatorially more training data

[illegible]

old  
form