IMARKUS LUDWIG I DATA COUNCIL I BARCELONA 2019 IMPROVING SEARCH WITH NATURAL LANGUAGE PROCESSING AND DEEP LEARNING





AUTO SCOUT 24	Suchen ∨	Ve	rkaufen ∨	Informier	en 🗸
	Martine .				
	~~~~				
✓ Neu ✓	Gebraucht	<ul> <li>.</li> </ul>	Tageszulassun	g	
Marke		~	Modell		~
• Verkaufsp	oreis 🔵 Rat	te	Preis bis (€)		~
Erstzulassun	g ab	~	Deutsch	nland	~
Stadt/PLZ	Umkreis	~	Grenzübe	ergreifend	
Weitere Suchoptionen           NEU         Smarte Textsuche			1.114.935 Treffer		

#### Finanzieren 🗡

Werkstatt

Mein Konto ∨ Jetzt anmelden **—** ~

Х

### Film ab für coole Auto-Tests!

Laufend neue Videos im Magazin.

Zu den Videos

Vite



#### Finanzieren 🗡

Werkstatt

Mein Konto ∨ Jetzt anmelden **—** ~

×

#### Film ab für coole Auto-Tests!

Laufend neue Videos im Magazin.

Zu den Videos

Vite

## Audi RS6 Performance Panorama 2018



make

## Audi RS6 Performance Panorama 2018



model

## Audi RS6 Performance Panorama 2018

make



model

version

## Audi RS6 Performance Panorama 2018

make





model

version

## Audi RS6 Performance Panorama 2018

make





model

version

## Audi RS6 Performance Panorama 2018

make

equipment

## first registration

filter

## Audi RS6 Performance Panorama 2018

filter

keyword







filter

# SEQUENCE-TO-SEQUENCE LEARNING

INPUT étudiant suis Je



Source: <u>http://jalammar.github.io/illustrated-transformer/</u>



# SEQUENCE-TO-SEQUENCE LEARNING

## natural language



Source: <a href="http://jalammar.github.io/illustrated-transformer/">http://jalammar.github.io/illustrated-transformer/</a>















peogot peogout peuceot peugeaut peugeot peugeu peugoet peugot peujeot peuscho

# VARIETY IN NATURAL LANGUAGE

peogot	rer
peogout	rer
peuceot	rer
peugeaut	rer
peugeot	rer
peugeu	rer
peugoet	rer
peugot	rer
peujeot	rer
peuscho	rin

# VARIETY IN NATURAL LANGUAGE

- nalt
- naud .
- naul
- nauld
- nault
- naults
- naut
- nolte
- noult
- nault

peogot	rer
peogout	rer
peuceot	rer
peugeaut	rer
peugeot	rer
peugeu	rer
peugoet	rer
peugot	rer
peujeot	rer
peuscho	rin

# VARIETY IN NATURAL LANGUAGE

nalt naud naul nauld

nault

naults

naut

nolte

noult

nault

volkswagem volkswagen wolcvagen wolksvagen wolsfagen wolsvagan

volcwagen volfsvagen volksagen volkswagem



## "In God we trust, all others bring data."

-W. EDWARDS DEMING

# TRAINING



e Pricing $\sim$	Search		Sign in	Sign up
	• Watch 437	★ Star 8,8	67 <b>%</b> Fork	c 2,288
Insights				
ake deep learning r nent-learning tpu	more accessible and a	ccelerate ML	research.	
S 67 releases	Le 192 contributo	ors	ঠাুঁ Apache-2	2.0
		Find File	Clone or de	ownload <del>-</del>
putations (both diagona	l and full)	Latest comm	it 9dff225 9 h	nours ago
			la	st month
py from oss_tests, th	is was deleted rece		la	st month
computations (both d	liagonal and full)		Qh	

# arXiv:1706.03762v5 [cs.CL] 6 Dec 2017

#### **Attention Is All You Need**

Ashish Vaswani* Google Brain avaswani@google.com Noam Shazeer* Google Brain

noam@google.com nikip@google.com usz@google.com

Llion Jones* Google Research llion@google.com Aidan N. Gomez^{*†} University of Toronto aidan@cs.toronto.edu

cs.toronto.edu lukaszkaiser@google.com
Illia Polosukhin* ‡

Niki Parmar*

Google Research

Jakob Uszkoreit*

**Google Research** 

Łukasz Kaiser*

Google Brain

**Oriol Vinyals** 

vinyals@google.com

DeepMind

#### Abstract

illia.polosukhin@gmail.com

The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show these models to be superior in quality while being more parallelizable and requiring significantly

Published as a conference paper at ICLR 2019

#### UNIVERSAL TRANSFORMERS

Mostafa Dehghani*†Stephan GUniversity of AmsterdamDeepMinddehghani@uva.nlsgouws@c

Stephan Gouws* n DeepMind sgouws@google.com

Jakob Uszkoreit Google Brain usz@google.com **Łukasz Kaiser** Google Brain lukaszkaiser@google.com

#### ABSTRACT

Recurrent neural networks (RNNs) sequentially process data by updating their state with each new data point, and have long been the de facto choice for sequence modeling tasks. However, their inherently sequential computation makes them slow to train. Feed-forward and convolutional architectures have recently been shown to achieve superior results on some sequence modeling tasks such as machine translation, with the added advantage that they concurrently process all inputs in the sequence, leading to easy parallelization and faster training times. Despite these successes, however, popular feed-forward sequence models like the Transformer fail to generalize in many simple tasks that recurrent models handle with ease, e.g. copying strings or even simple logical inference when the string or formula lengths exceed those observed at training time. We propose the Universal Transformer (UT), a parallel-in-time self-attentive recurrent sequence model which can be cast as a generalization of the Transformer model and which addresses these issues. UTs combine the parallelizability and global receptive field of feed-forward sequence models like the Transformer with the recurrent inductive bias of RNNs. We also add a dynamic per-position halting mechanism and find that it improves accuracy on several tasks. In contrast to the standard Transformer, under certain assumptions UTs can be shown to be Turing-complete. Our experiments show that UTs outperform standard Transformers on a wide range of algorithmic and language understanding tasks, including the challenging LAMBADA language modeling task where UTs achieve a new state of the art, and machine translation where UTs achieve a 0.9 BLEU improvement over Transformers on the WMT14 En-De dataset.

1 INTRODUCTION

iv:1807.03819v3 [cs.CL] 5 Mar 2019

#### **Fast Decoding in Sequence Models Using Discrete Latent Variables**

Łukasz Kaiser¹ Aurko Roy¹ Ashish Vaswani¹ Niki Parmar¹ Samy Bengio¹ Jakob Uszkoreit¹ Noam Shazeer¹

#### Abstract

Autoregressive sequence models based on deep neural networks, such as RNNs, Wavenet and the Transformer attain state-of-the-art results on many tasks. However, they are difficult to parallelize and are thus slow at processing long sequences. RNNs lack parallelism both during training and decoding, while architectures like WaveNet and Transformer are much more parallelizable during training, yet still operate sequentially during decoding. We present a method to extend sequence models using discrete latent variables that makes decoding much more parallelizable. We first autoencode the target sequence into a shorter sequence of discrete latent variables, which at inference time is generated autoregressively, and finally decode the output sequence from this shorter latent sequence in parallel. To this end, we introduce a novel method for constructing a sequence of discrete latent variables and compare it with previously introduced methods. Finally, we evaluate our model end-to-end on the task of neulation (Sutskever et al., 2014; Bahdanau et al., 2014; Cho et al., 2014), parsing (Vinyals et al., 2015), and many others. RNNs are inherently sequential, however, and thus tend to be slow to execute on modern hardware optimized for parallel execution. Recently, a number of more parallelizable sequence models were proposed and architectures such as WaveNet (van den Oord et al., 2016), ByteNet (Kalchbrenner et al., 2016) and the Transformer (Vaswani et al., 2017) can indeed be trained faster due to improved parallelism.

When actually generating sequential output, however, their autoregressive nature still fundamentally prevents these models from taking full advantage of parallel computation. When generating a sequence  $y_1 \ldots y_n$  in a canonical order, say from left to right, predicting the symbol  $y_t$  first requires generating all symbols  $y_1 \ldots y_{t-1}$  as the model predicts

 $P(y_t|y_{t-1} y_{t-2} \dots y_1).$ 

During training, the ground truth is known so the conditioning on previous symbols can be parallelized. But during decoding, this is a fundamental limitation as at least n sequential steps need to be made to generate  $y_1 \dots y_n$ .

#### **The Evolved Transformer**

**David R. So**¹ Chen Liang¹ Quoc V. Le¹

#### Abstract

Recent works have highlighted the strength of the Transformer architecture on sequence tasks while, at the same time, neural architecture search (NAS) has begun to outperform human-designed models. Our goal is to apply NAS to search for a better alternative to the Transformer. We first construct a large search space inspired by the recent advances in feed-forward sequence models and then run evolutionary architecture search with warm starting by seeding our initial population with the Transformer. To directly search on the computationally expensive WMT 2014 English-German translation task, we develop the Progressive Dynamic Hurdles method, which allows us to dynamically allocate more resources to more promising candidate models. The architecture found in our experiments - the Evolved Transformer – demonstrates consistent improvement over the Transformer on four well-established language tasks: WMT 2014 English-German, WMT 2014 English-French, WMT 2014 English-Czech and LM1B. At a big model size, the Evolved Transformer establishes a new state-ofthe-art BLEU score of 29.8 on WMT'14 English-German; at smaller sizes, it achieves the same quality as the original "big" Transformer with 37.6% less parameters and outperforms the Transformer by 0.7 BLEU at a mobile-friendly model

models, although some effort has also been invested in searching for sequence models (Zoph & Le, 2017; Pham et al., 2018). In these cases, it has always been to find improved recurrent neural networks (RNNs), which were long established as the de facto neural model for sequence problems (Sutskever et al., 2014; Bahdanau et al., 2015).

However, recent works have shown that there are better alternatives to RNNs for solving sequence problems. Due to the success of convolution-based networks, such as Convolution Seq2Seq (Gehring et al., 2017), and full attention networks, such as the Transformer (Vaswani et al., 2017), feed-forward networks are now a viable option for solving sequence-to-sequence (seq2seq) tasks. The main strength of feed-forward networks is that they are faster, and easier to train than RNNs.

The goal of this work is to examine the use of neural architecture search methods to design better feed-forward architectures for seq2seq tasks. Specifically, we apply *tournament selection* architecture search and *warm start* it with the Transformer, considered to be the state-of-art and widely-used, to evolve a better and more efficient architecture. To achieve this, we construct a search space that reflects the recent advances in feed-forward seq2seq models and develop a method called *Progressive Dynamic Hurdles* (PDH) that allows us to perform our search directly on the computationally demanding WMT 2014 English-German (En-De) translation task. Our search produces a new architecture – called the *Evolved Transformer* (ET) –

6

201

May

 $\sim$ 

5

[cs.L

1117v4

:1901.1

Lean

### 

pip install tensor2tensor

t2t-trainer \
 --generate_data \
 --model=\$MODEL \
 --problem=\$PROBLEM \
 --data_dir=\$DATA_DIR \
 --hparams_set=\$HPARAMS \
 --output_dir=\$TRAIN_DIR





# DEPLOYMENT API

	Why GitHub? $\sim$	Enterprise	Explore $\sim$	Marketplace
📮 enc <> Co	ode / <b>starlette</b> de ① Issues 46	ື່ 1) Pull requ	uests 39	Projects 0
The litt python	le ASGI framewor	k that shines ckets graphq	. 💥 https	://www.starlet
	<b>592</b> commits	<b>پ</b> 18	5 branches	
Branch	n: master 👻 New pu	Ill request		
to	mchristie Merge pull re	equest #620 from	elyobo/docum	nent-request-app-
🖬 do	CS	Merge	pull request	#620 from elyo
scr	ipts	Reinst	ate Lifespan	Middleware, for
🖬 sta	rlette	Versio	n <b>0.12.9</b>	



# DEPLOYMENT MODEL



Source: <u>https://www.tensorflow.org/tfx/serving/architecture</u>



# THANK YOU FOR YOUR ATTENTION!

https://www.linkedin.com/in/markusludwig/