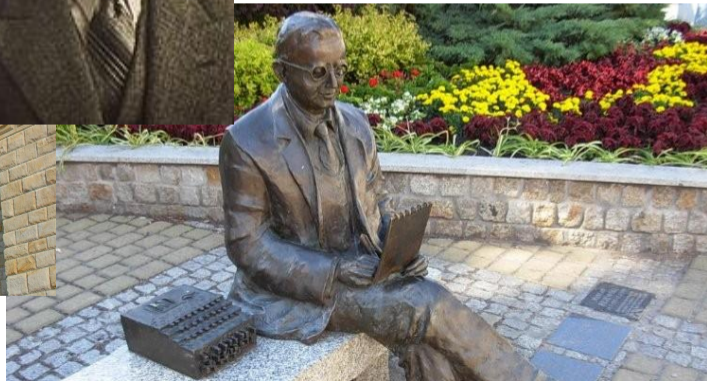
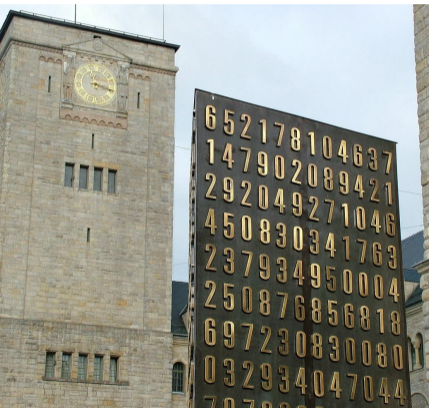


Marian: Homecoming

Marcin Junczys-Dowmunt

Translator

Marian Rejewski



A few words about Marian

- Portable C++ code with minimal dependencies (CUDA or MKL and still Boost);
- Single engine for training and decoding on GPU and CPU;
- Custom auto-diff engine with dynamic graphs (similar to DyNet);
- Optimized towards NMT.
- `http://marian-nmt.github.io` and
`https://github.com/marian-nmt/marian`

Part I

A Machine Translation Marathon 2016 project

The first commit

Commit: 6a7c93

Date: May 4th, 2016

Message: very cool

Lines: 155

```
#include <iostream>
#include "mad.h"

int main(int argc, char** argv) {
    Var x0 = 1, x1 = 2, x2 = 3;
    auto y = x0 + x0 + log(x2) + x1;

    std::vector<Var> x = { x0, x1, x2 };

    set_zero_all_adjoints();
    y.grad();

    std::cerr << "y_=" << y.val() << std::endl;
    for(int i = 0; i < x.size(); ++i)
        std::cerr << "dy/dx_" << i << "_="
            << x[i].adj() << std::endl;
}
```

```
Var x0 = 1, x1 = 2, x2 = 3;  
auto y = x0 + x0 + log(x2) + x1;
```

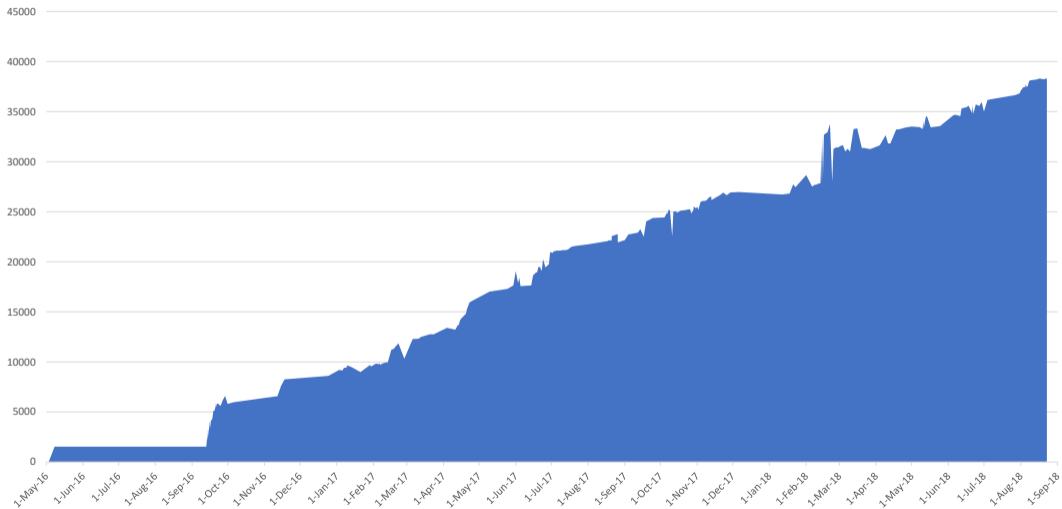
```
y = 5.09861
```

```
dy/dx_0 = 2
```

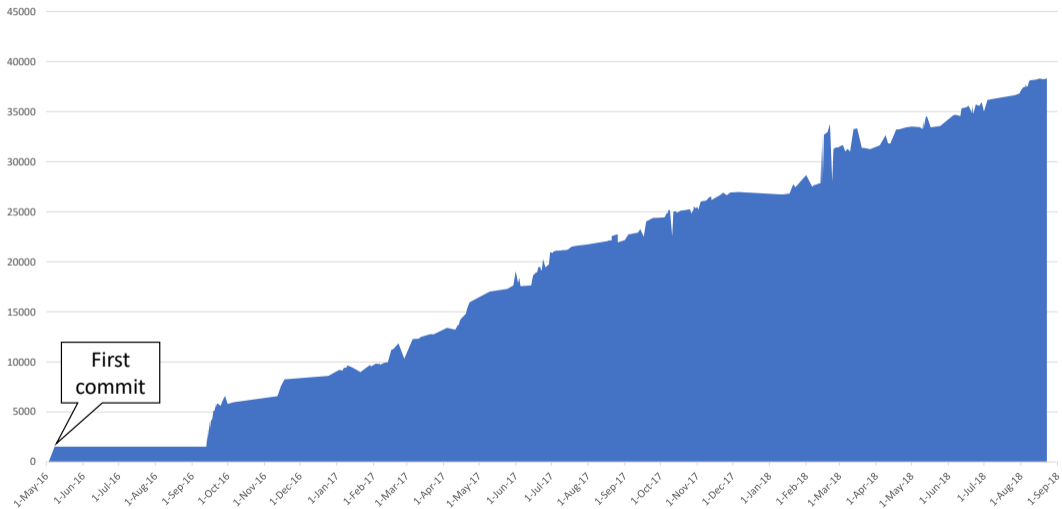
```
dy/dx_1 = 1
```

```
dy/dx_2 = 0.333333
```

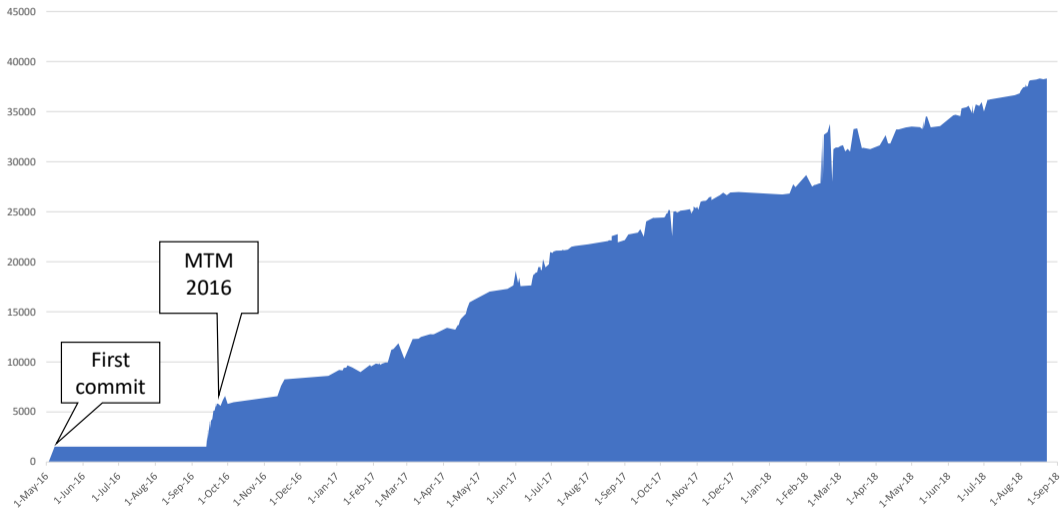
Lines of code over time



Lines of code over time



Lines of code over time



A Neural Network Toolkit for MT

Maximiliana Behnke
Tomasz Dwojak
Marcin Junczys-Dowmunt
Roman Grundkiewicz
Andre Martins
Hieu Hoang
Lane Schwartz

A Neural Network Toolkit ~~for MT~~

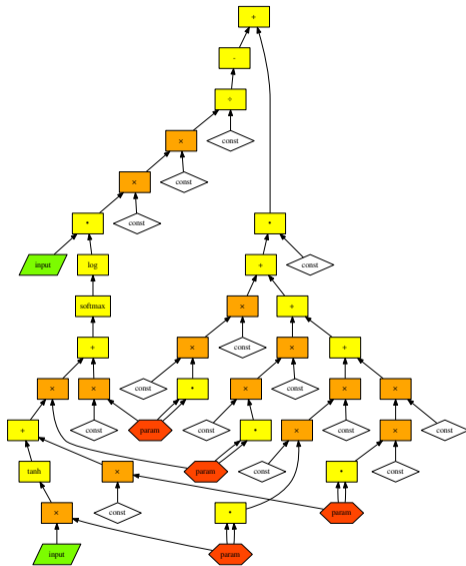
Maximiliana Behnke
Tomasz Dwojak
Marcin Junczys-Dowmunt
Roman Grundkiewicz
Andre Martins
Hieu Hoang
Lane Schwartz

Why create another NN toolkit?

- Flexibility
 - Add functionality easier & faster
- Speed
 - Pure C++ implementation
 - GPU-enabled (CPU may come soon)
- Learn about Deep Learning
 - Implement everything from scratch by ourselves

What we've achieved this week?

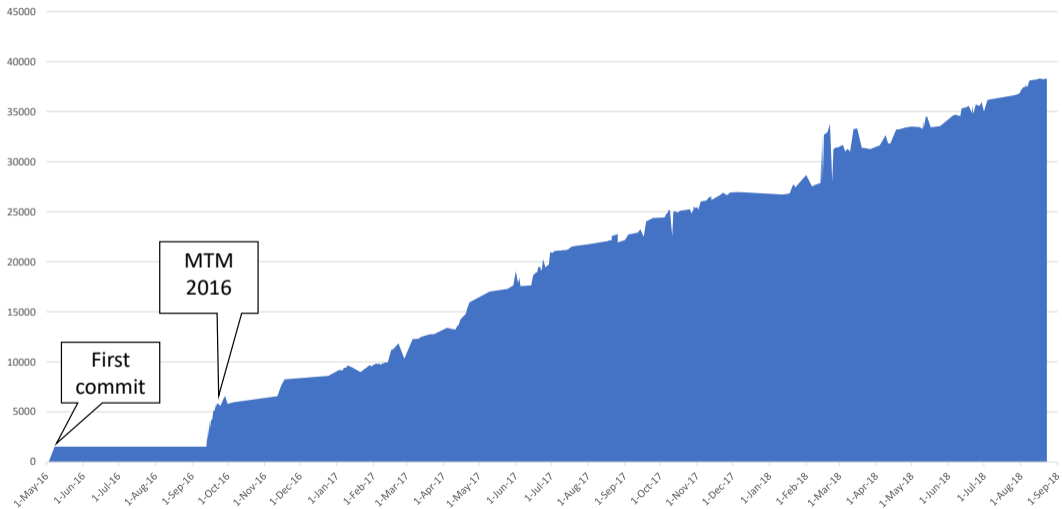
- Framework to create computation graphs
 - Simple feedforward NN
 - RNN, GRU, LSTM...
 - Binary or multiclass classifier
- Forward step
 - Classify, given input data and weights
- Backward step
 - Learn weights using backpropagation
- Tested with small datasets
 - MNIST (digit image recognition task)
 - MT
- Documentation
 - Doxygen



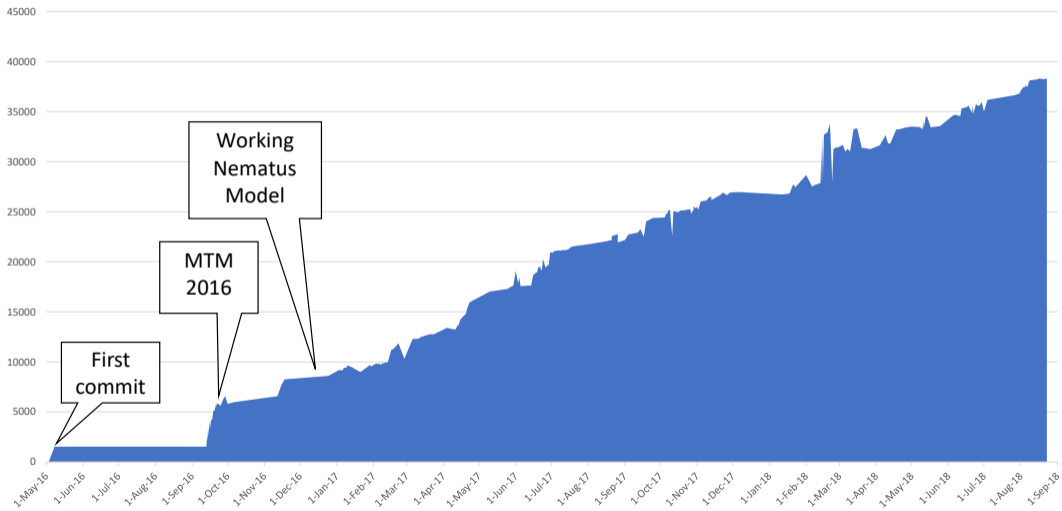
What needs to be finished?

- Basic features:
 - Data shuffling
 - More random distributions
 - ...
- Model serialization & deserialization
- Documentation
- ...

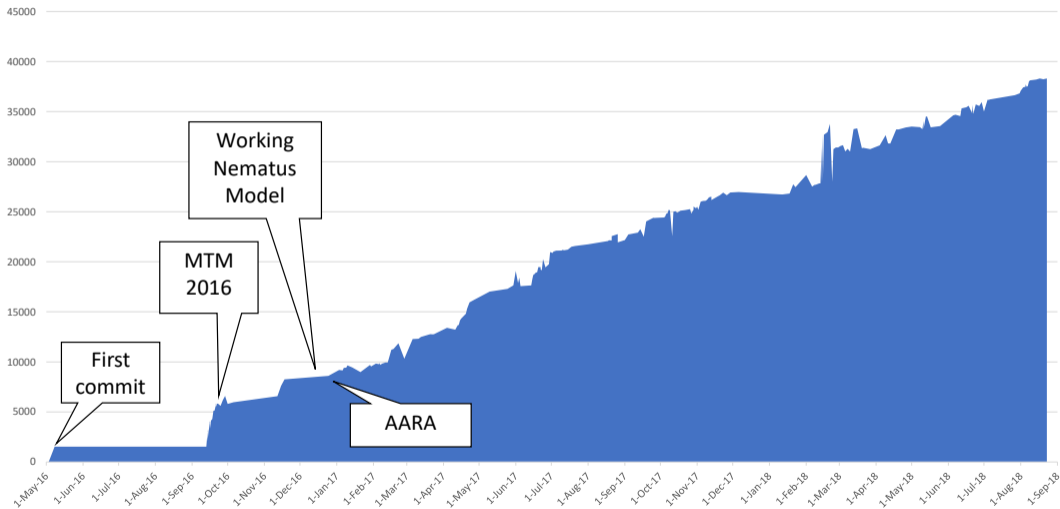
Lines of code over time



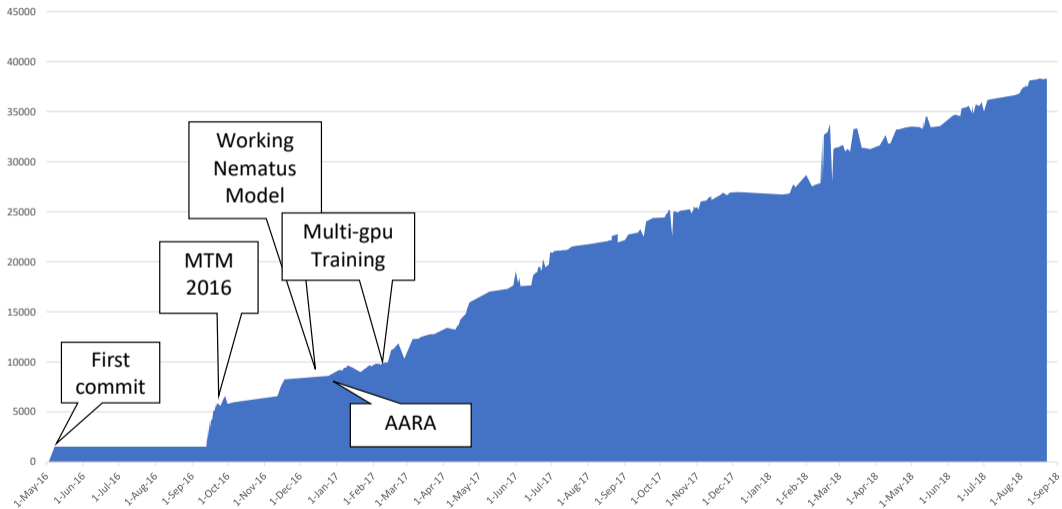
Lines of code over time



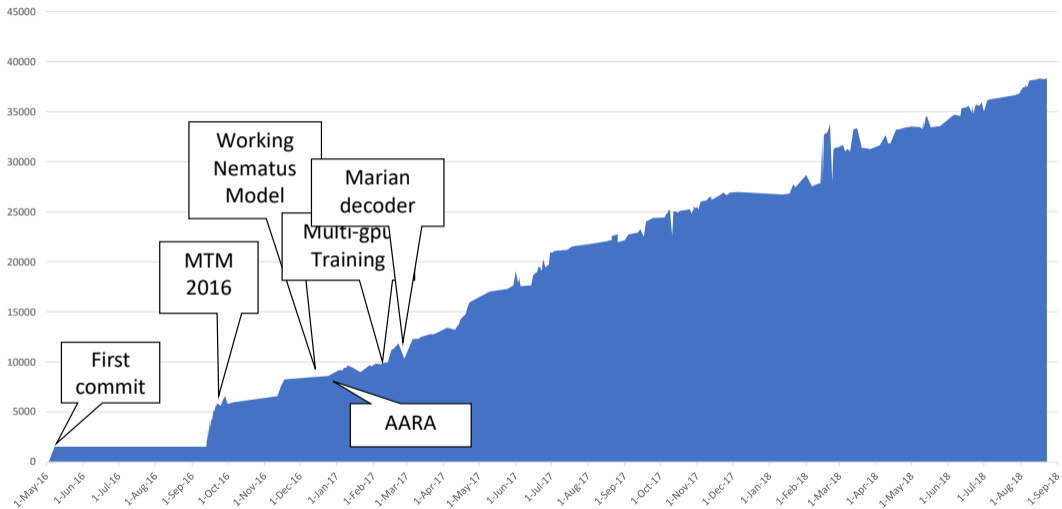
Lines of code over time



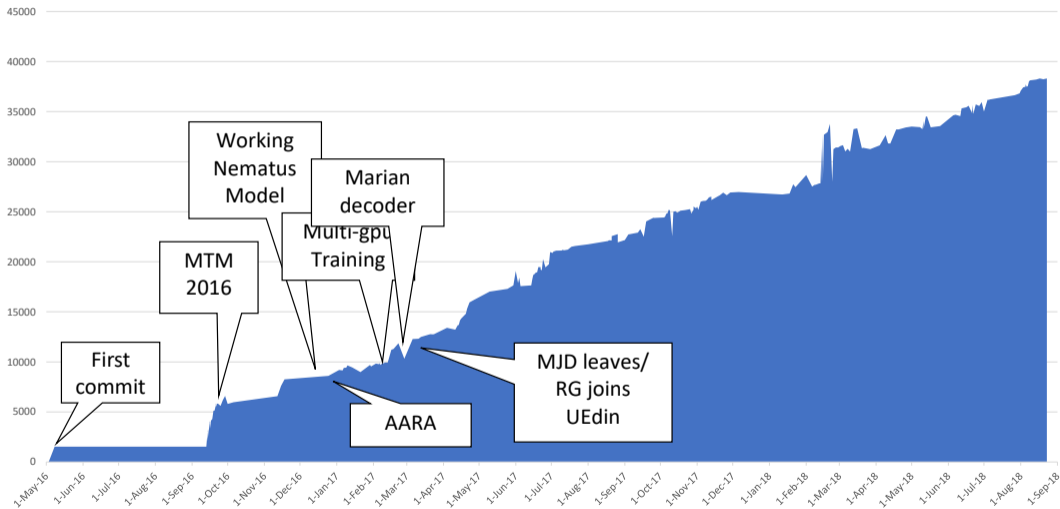
Lines of code over time



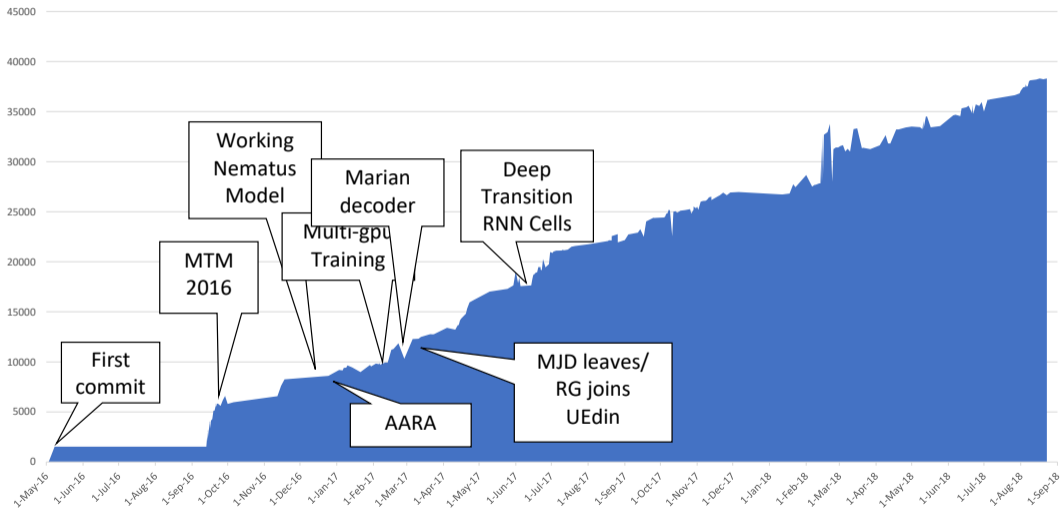
Lines of code over time



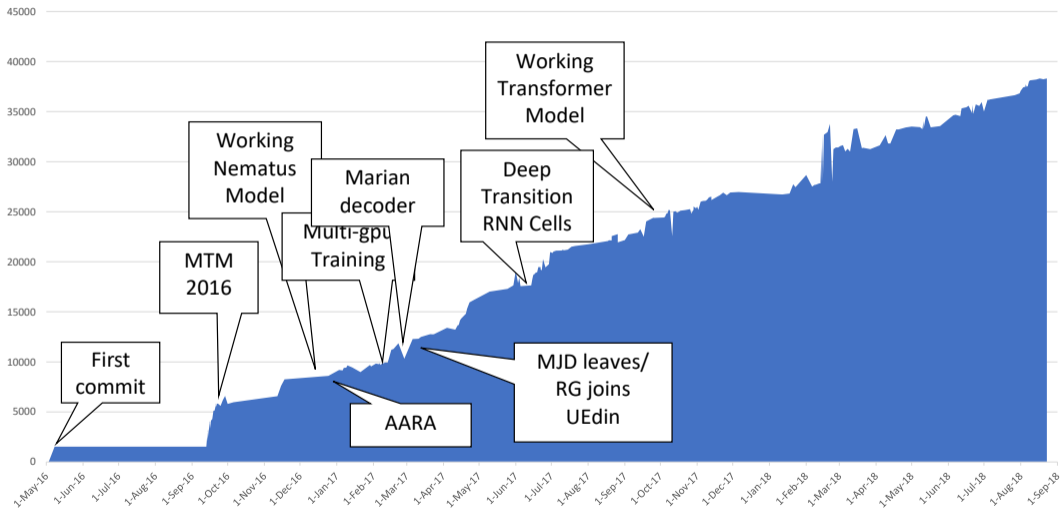
Lines of code over time



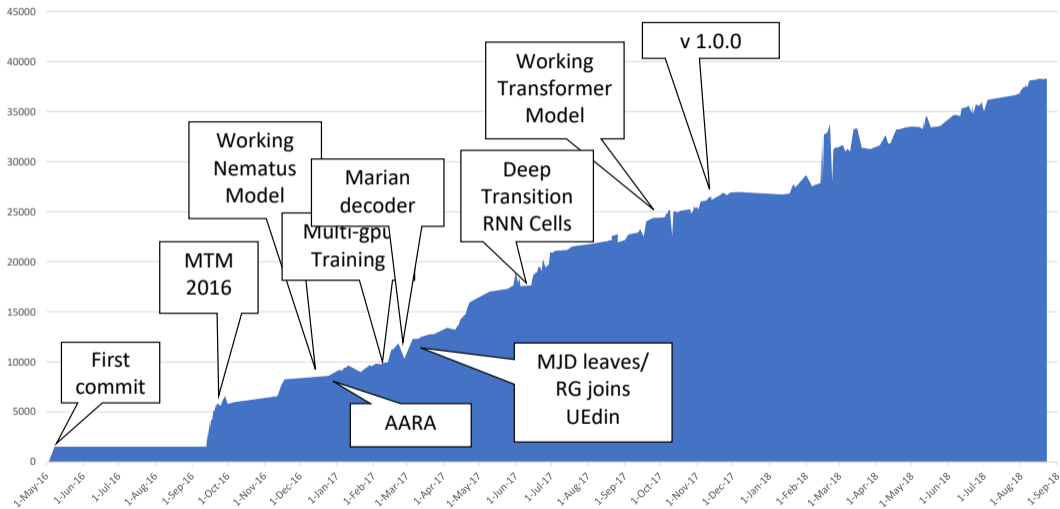
Lines of code over time



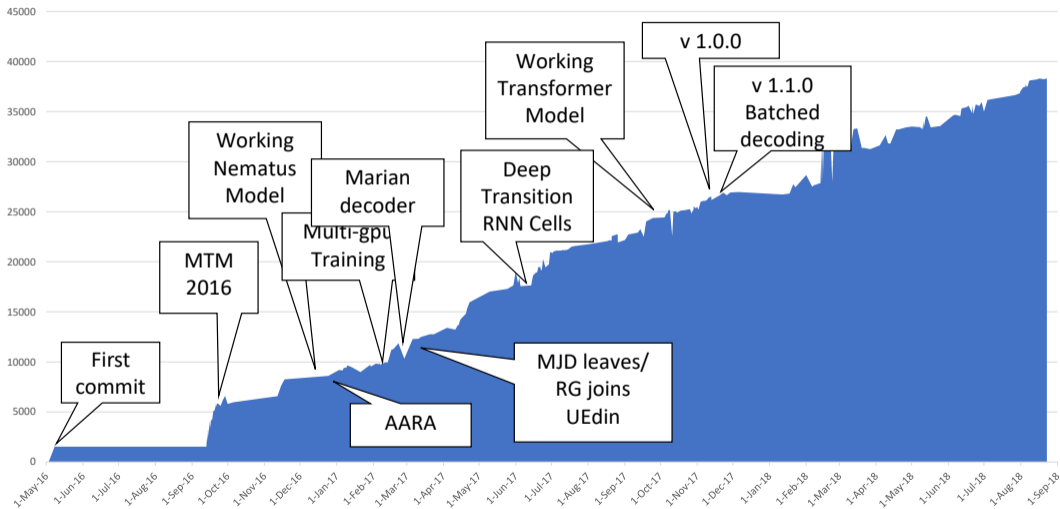
Lines of code over time



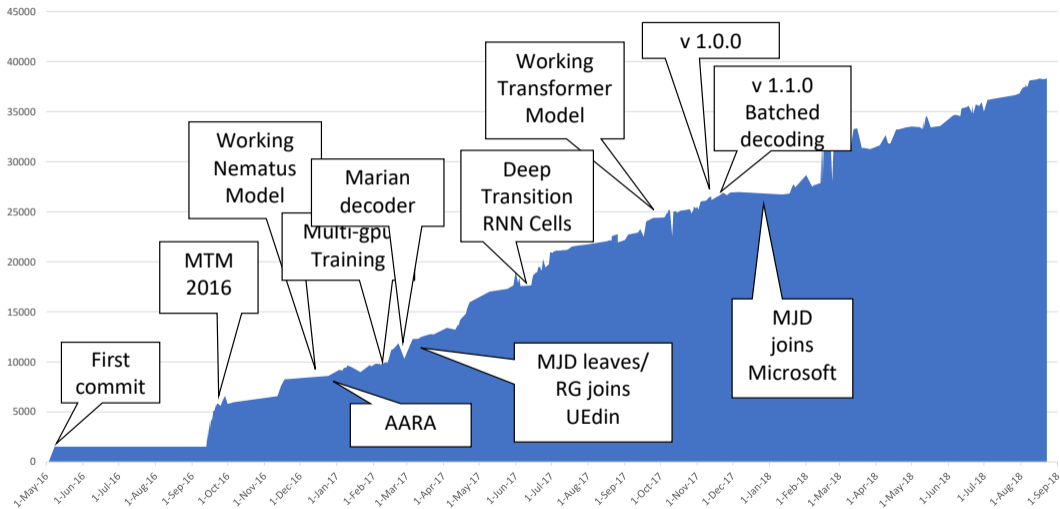
Lines of code over time



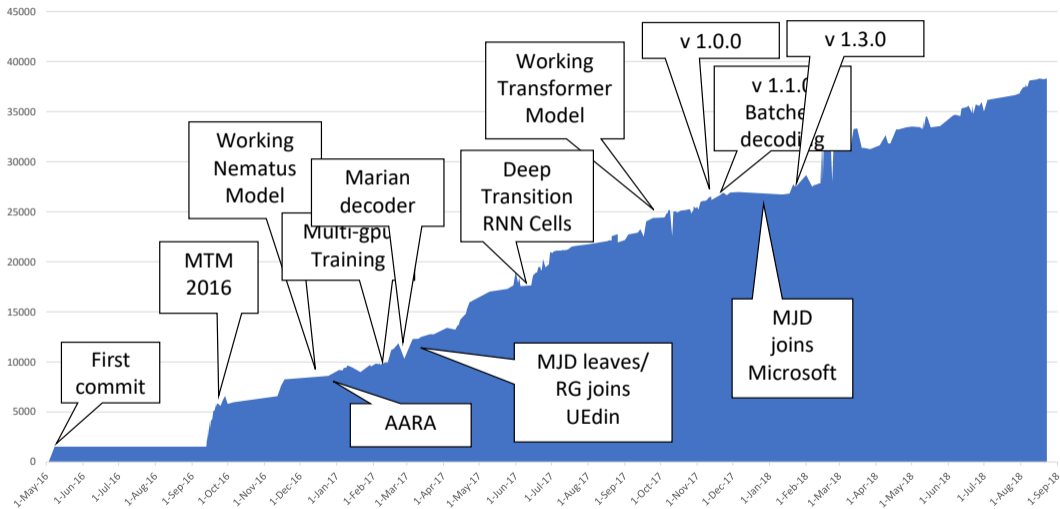
Lines of code over time



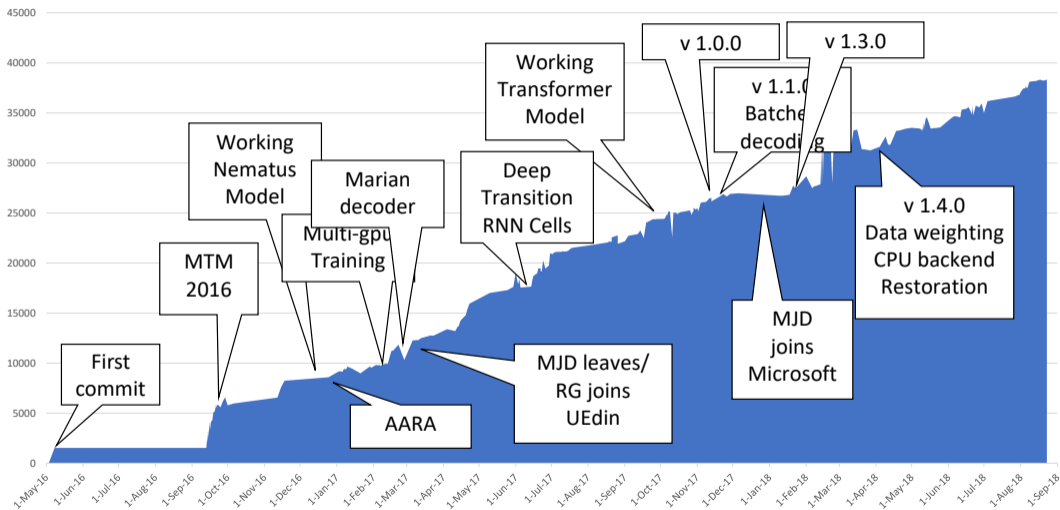
Lines of code over time



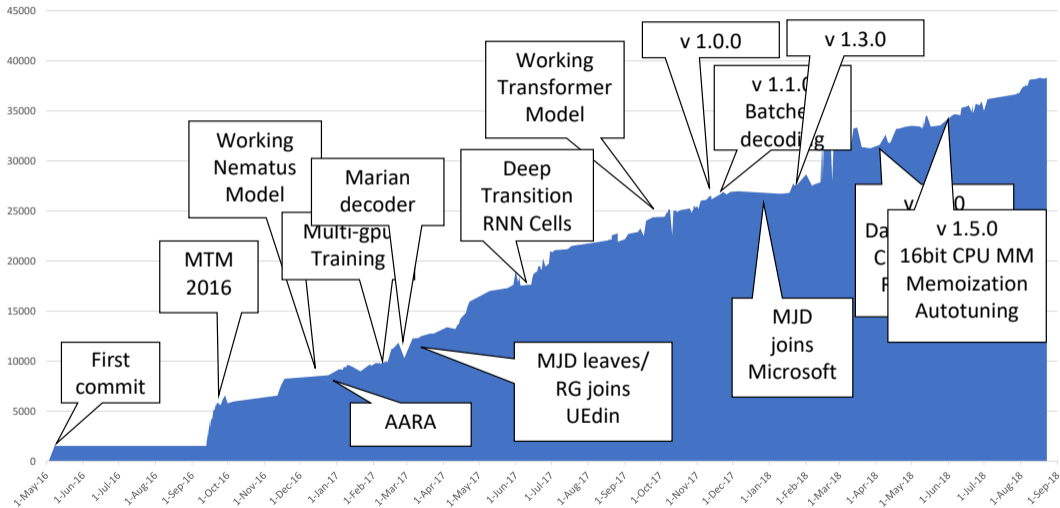
Lines of code over time



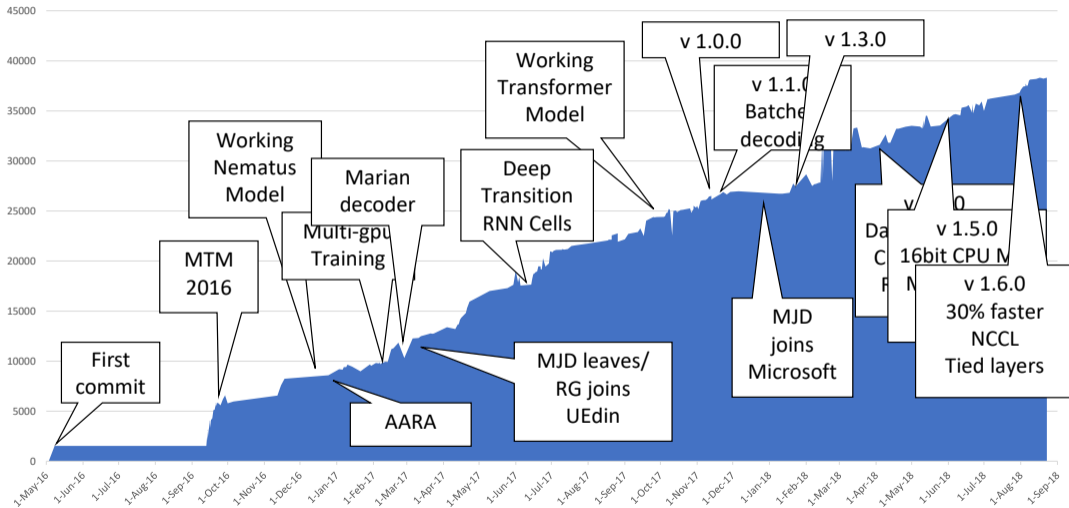
Lines of code over time



Lines of code over time

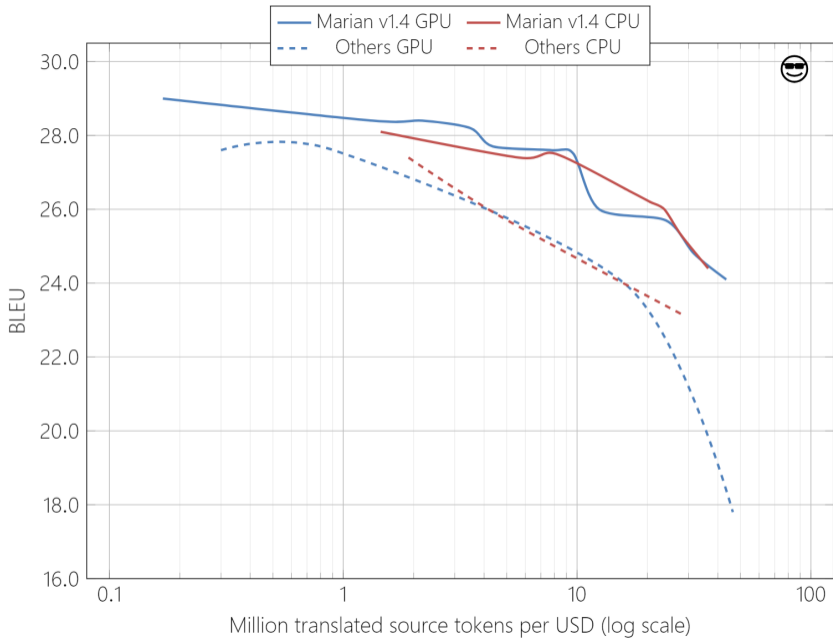


Lines of code over time



Going further

- Reduce dependencies for CPU version to zero
- Reduce dependencies for GPU version to CUDA
- Become faster and more versatile
- Research tool with immediate deployment





System Output List

English-German newstest2018

[Translations](#) [Resources](#) [Download](#) [Info](#) [Account](#)

Interested in Contributing?

- Check out available [resources](#).
- Create an [account](#) and start [submitting](#) your own systems.

Scored Systems

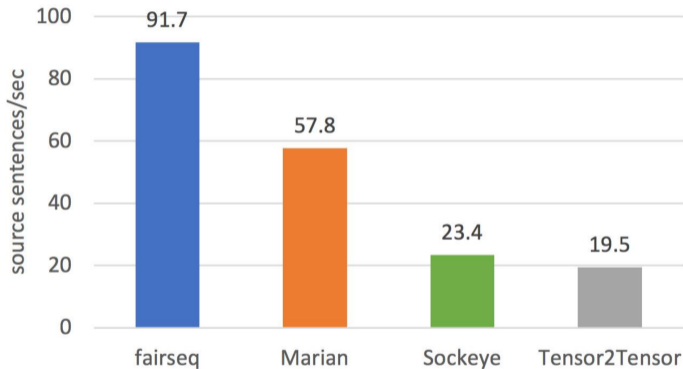
System	Submitter	System Notes	Constraint	Run Notes	BLEU	BLEU-cased	TER	BEER 2.0	CharactER
Marian-Transformer (Details)	marcinjd Microsoft	Marian Transformer-Big	yes	Transformer-big ensemble x4. With back-translation, data-filtering on Paracrawl with domain-weighting. Decoder-time ensembling with transformer-LM, right-to-left decoding.	48.9	48.3	0.407	0.697	0.362
NMT-SMT_Hybrid (Details)	fstahlberg University of Cambridge	MBR-based combination of neural models and SMT	yes		47.1	46.6	0.415	0.691	0.369
NTT Transformer-based System (Details)	makoto-mr NTT	Based on Transformer Big model. Trained with filtered version of CommonCrawl, ParaCrawl and synthetic corpus of newscrawl2017. R2L reranking.	yes		47.0	46.5	0.426	0.688	0.370
Primary Submission (Details)	pianist Karlsruhe Institute of Technology	Primary Submission	yes		46.9	46.3	0.428	0.687	0.382
MMT-unconstraint (Details)	nicolabertoldi MMT srl		no		46.7	46.2	0.432	0.682	0.387
JHU (Details)	jhu-nmt Johns Hopkins University	Marian Deep RNN	yes	Marian deep model, ensemble of 4 runs using base data (without Paracrawl), re-back-translated news 2016. Not final system yet.	43.6	43.0	0.453	0.670	0.394
MMT-constraint-dec (Details)	nicolabertoldi MMT srl		yes		42.9	42.5	0.463	0.667	0.411
NJUNMT (Details)	ZhaoChengqi Nanjing University	transformer base without back translation	yes	transformer base without back translation	40.6	40.0	0.496	0.647	0.436
LMU-unsupervised-nmt-wmt18-en-de (Details)	Matthias Huck LMU Munich	Unsupervised NMT (no parallel training corpora)	yes		15.8	15.5	0.762	0.500	failed
RWTH Unsupervised NMT Ensemble (Details)	yunsukim RWTH Aachen University	(Unsupervised) Transformer with shared encoder/decoder, separate top-50k word	yes		15.9	14.8	0.753	0.514	0.607

Interested in Contributing?

- Check out available [resources](#).
- Create an [account](#) and start [submitting](#) your own systems.

Scored Systems

System	Submitter	System Notes	Constraint	Run Notes	BLEU	BLEU-cased	TER	BEER 2.0	CharactER
Marian-Transformer (Details)	marcinjd Microsoft	Marian Transformer-Big	yes	Transformer-big ensemble x4. With back-translation, data-filtering on Paracrawl with domain-weighting Decoder-time ensembling with transformer-LM, right-to-left decoding.	48.9	48.3	0.407	0.697	0.362
NMT-SMT-Hybrid (Details)	fstahlberg University of Cambridge	MBR-based combination of neural models and SMT	yes		47.1	46.6	0.415	0.691	0.369
NTT Transformer-based System (Details)	makoto-mr NTT	Based on Transformer Big model. Trained with filtered version of CommonCrawl, ParaCrawl and synthetic corpus of newscrawl2017. R2L reranking.	yes		47.0	46.5	0.426	0.688	0.370
Primary Submission (Details)	pianist Karlsruhe Institute of Technology	Primary Submission	yes		46.9	46.3	0.428	0.687	0.382
MMT-unconstrained (Details)	nicolabertoldi MMT srl		no		46.7	46.2	0.432	0.682	0.387
JHU (Details)	jhu-nmt Johns Hopkins University	Marian Deep RNN	yes	Marian deep model, ensemble of 4 runs using base data (without Paracrawl), re-back-translated news 2016. Not final system yet.	43.6	43.0	0.453	0.670	0.394
MMT-constraint-dec (Details)	nicolabertoldi MMT srl		yes		42.9	42.5	0.463	0.667	0.411
NJUNMT (Details)	ZhaoChengqi Nanjing University	transformer base without back translation	yes	transformer base without back translation	40.6	40.0	0.496	0.647	0.436
LMU-unsupervised-nmt-wmt18-en-de (Details)	Matthias Huck LMU Munich	Unsupervised NMT (no parallel training corpora)	yes		15.8	15.5	0.762	0.500	failed
RWTH Unsupervised NMT Ensemble (Details)	yunsukim RWTH Aachen University	(Unsupervised) Transformer with shared encoder/decoder, separate top-50k word	yes		15.9	14.8	0.753	0.514	0.607



Michael Auli

June 15 · 🌐

We are releasing new features for fairseq, FAIR's sequence to sequence learning library: <https://github.com/pytorch/fairseq>

Distributed training, fp16, delayed batching
 We release code and pre-trained models to reproduce our recent paper "Scaling Neural Machine Translation" (<https://arxiv.org/abs/1806.00187>) where we train on up to 128 GPUs with half precision floating point operations as well ... [See More](#)

👍❤️👍 133

49 Shares



Like



Comment



Share



Write a comment...



Fast inference

Fairseq can generate translations at a rate of 92 sentences/second for big Transformers on a fast GPU by clever caching, removing finished sentences from the computation and by batching by tokens. This improves speed by nearly 90%. The image shows a comparison to other stacks (measured on a V100 GPU for WMT English-German translation on newstest2014 using a big Transformer).

Language models

Fairseq now supports the training of gated convolutional language models (<https://arxiv.org/abs/1812.08020>). It can train a Google Billion Word language model on 128 GPUs in less than a day.



James Bradbury

@jekbradbury

Following



Facebook's fairseq MT engine is really, really fast... Like, 50% faster than [@marian_nmt](#) (which is itself way faster than Sockeye/OpenNMT/Tensor2Tensor/xnmt/Ne matus/etc) at generating from the same Transformer model
[facebook.com/61013326/posts ...](https://facebook.com/61013326/posts...)

2:24 PM - 15 Jun 2018

37 Retweets 139 Likes



2



37



139



Tweet your reply



Marian NMT @marian_nmt · Jun 15



Replying to [@jekbradbury](#)

Hold my beer ;)



1



13





James Bradbury

@jekbradbury

Following



Facebook's fairseq MT engine is really, really fast... Like, 50% faster than @marian_nmt (which is itself way faster than Sockeye/OpenNMT/Tensor2Tensor/xnmt/Ne matus/etc) at generating from the same Transformer model

[facebook.com/61013326/posts ...](https://facebook.com/61013326/posts...)

2:24 PM - 15 Jun 2018

37 Retweets 139 Likes



2



37



139



Tweet your reply



Marian NMT @marian_nmt · Jun 15

Replying to @jekbradbury

Hold my beer ;)



1



13



Marian NMT @marian_nmt · Jun 17

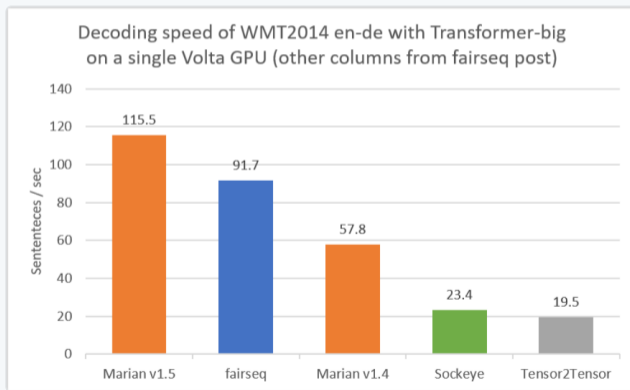
Boom! Marian v1.5.0 released.

Includes:

- Extensions from the WMT shared task on efficiency arxiv.org/abs/1805.12096
- Optimized GPU-decoding for Transformer models.

See chart below for speed comparison to v1.4.0 (based on FAIR's post)

[@jekbradbury](#) [@alvations](#)



4



33



108



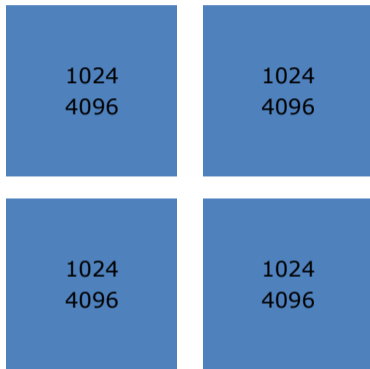
Part II

Decoding on the CPU

Quality first – speed later

- Lessons from WNMT shared task on efficient decoding;
- Sequence-level knowledge distillation (Kim & Rush 2016):
- Training four Transformer-big models on official task data (teacher);
- Translate entire EN data to DE-trans (8-best list);
- Select sentences with highest sentence-level BLEU based on DE-orig data;
- Train students on EN/DE-trans data.

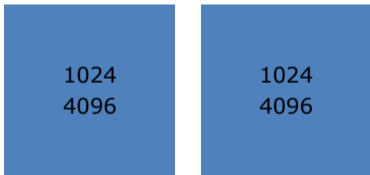
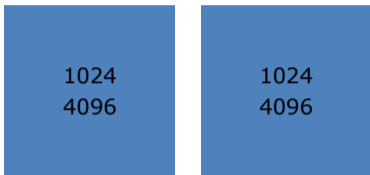
Teacher Transformer Big



4 × 813 MiB
29.0

(Scale preserving)

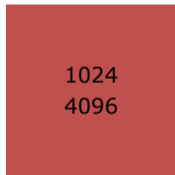
Teacher Transformer Big



4 × 813 MiB
29.0

Students Transformer beam=1

Big



813 MiB
28.2

Base



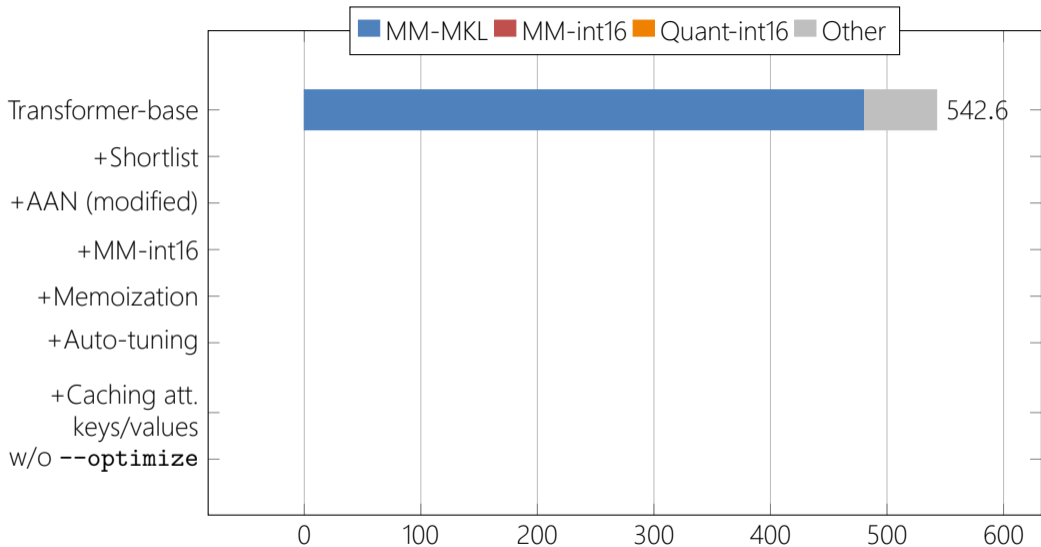
238 MiB
27.6

Small



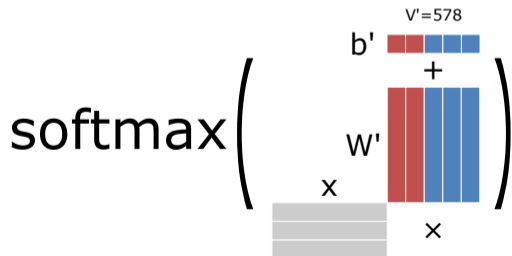
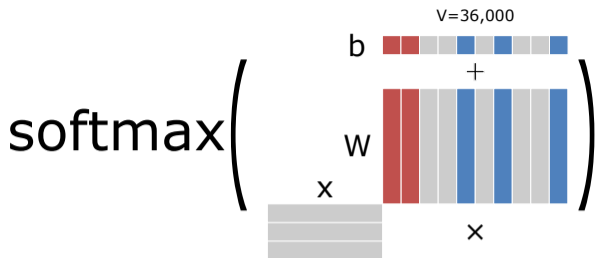
101 MiB
26.4

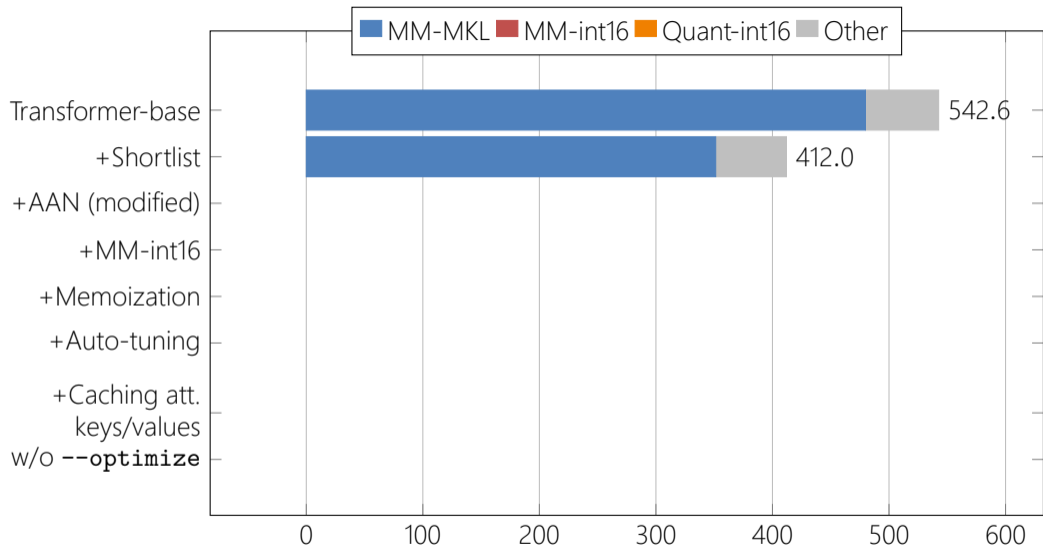
(Scale preserving)



Seconds to translate newstest2014 (batch-size: ca. 384 words – 5 to 25 sentences)

</s>	dog	Hund
<unk>	dog	Hun@@
,	dog	de@@
.	dog	Hunde
the	...	
in	lift	heben
die	lift	Li@@
of	lift	ft
der	lift	Aufzug
and	...	





Seconds to translate newstest2014 (batch-size: ca. 384 words – 5 to 25 sentences)

Multiplicative attention:

$$Q' = QW_q + b_q$$

$$K' = KW_k + b_k$$

$$V' = VW_v + b_v$$

$$C = \text{softmax}(Q' \times (K')^T) \times V'$$

$$Y = \text{norm}(Q + C)$$

- During training: $Q = K = V$
- During translation: $Q \neq K; K = V$
- Complexity per step: $O(n)$
- Because:

$$c_t = \text{softmax}(q'_t \times (K'_{<t})^T) \times V'_{<t}$$

$$K_{<t+1} = [K_{<t}; q_t]$$

Multiplicative attention:

$$Q' = QW_q + b_q$$

$$K' = KW_k + b_k$$

$$V' = VW_v + b_v$$

$$C = \text{softmax}(Q' \times (K')^T) \times V'$$

$$Y = \text{norm}(Q + C)$$

- During training: $Q = K = V$
- During translation: $Q \neq K; K = V$
- Complexity per step: $O(n)$
- Because:

$$c_t = \text{softmax}(q'_t \times (K'_{<t})^T) \times V'_{<t}$$

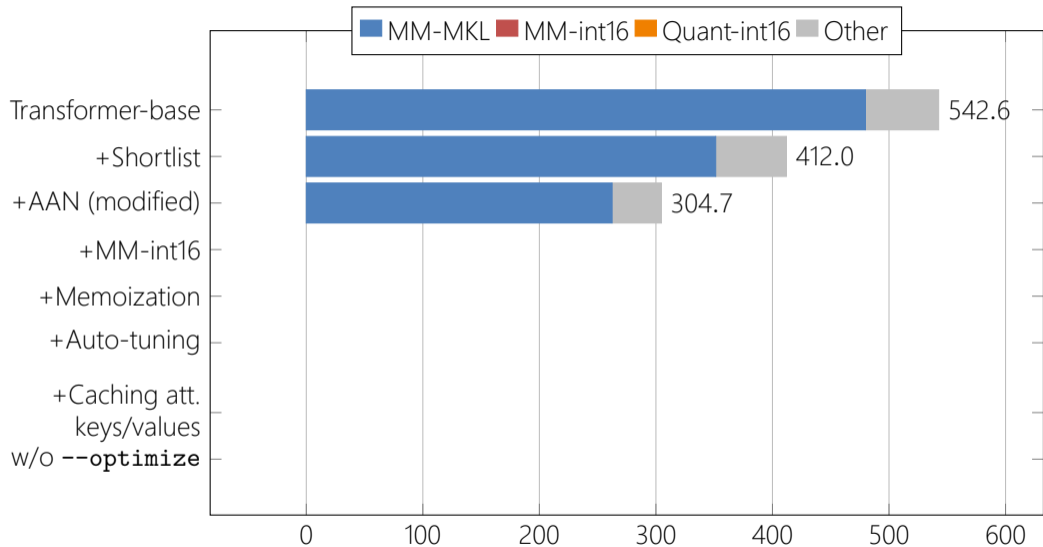
$$K_{<t+1} = [K_{<t}; q_t]$$

Average attention network
(Zhang et al. 2018):

$$C = \text{gate}(\text{FFN}(\bar{V}), Q)$$

$$Y = \text{norm}(Q + C)$$

- Gate and FFN optional
- Complexity per step: $O(1)$
- Because: $\bar{v}_t = \frac{1}{t}((t-1)\bar{v}_{t-1} + v_t)$
- Basically a weird RNN
- Authors report 4x speed-up for beam=4



Seconds to translate newstest2014 (batch-size: ca. 384 words – 5 to 25 sentences)

Code based on Devlin (2017), extended to AVX512

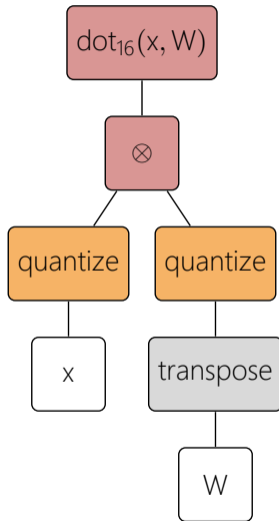
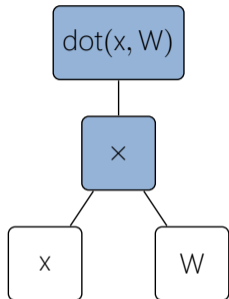
$$q(x) = \text{int16}(x \cdot 2^{10})$$

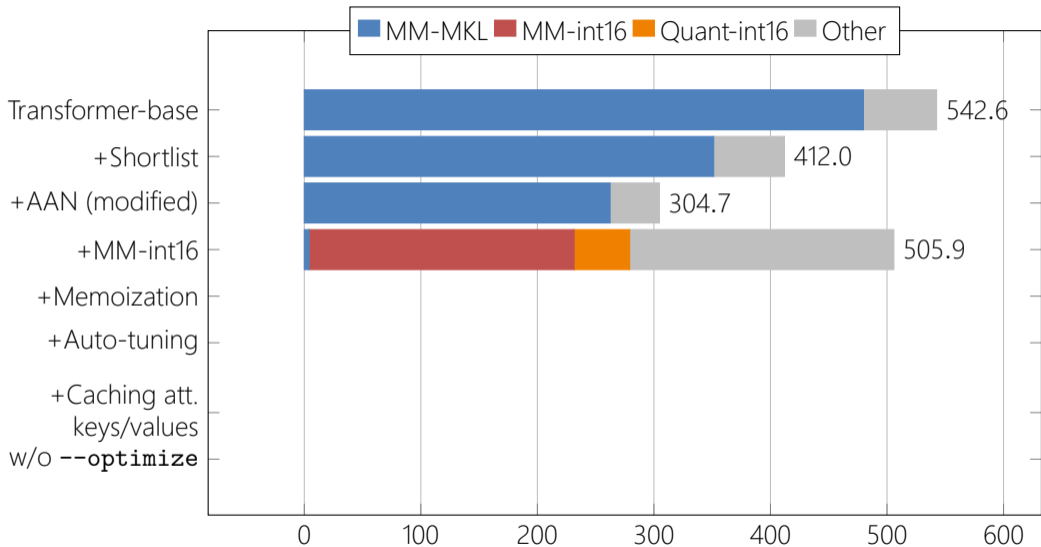
Code based on Devlin (2017), extended to AVX512

$$q(x) = \text{int16}(x \cdot 2^{10})$$
$$A_q \otimes B_q = A_q \times B_q^T$$

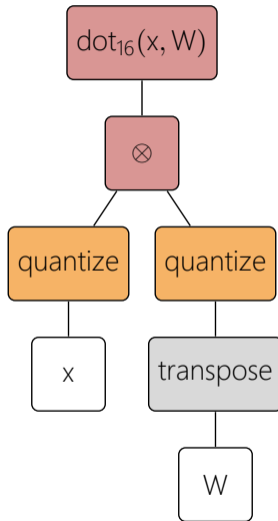
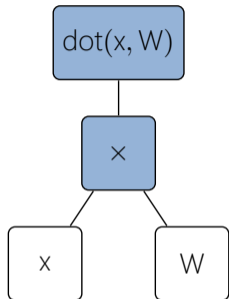
Code based on Devlin (2017), extended to AVX512

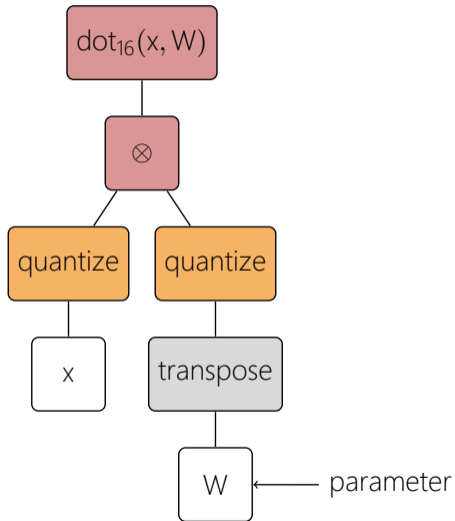
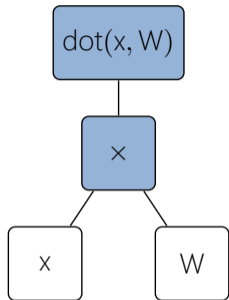
$$\begin{aligned}q(x) &= \text{int16}(x \cdot 2^{10}) \\A_q \otimes B_q &= A_q \times B_q^T \\ \\x \times W &= x \times (W^T)^T \\ &\approx q(x) \times q(W^T)^T \\ &= q(x) \otimes q(W^T)\end{aligned}$$

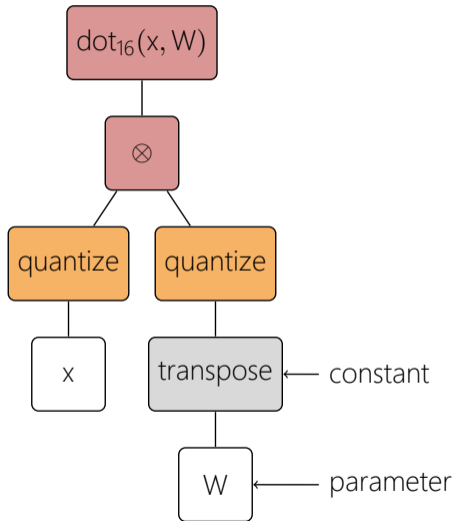
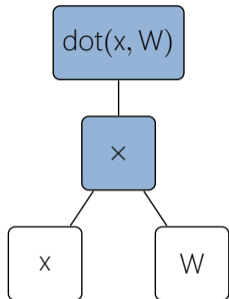


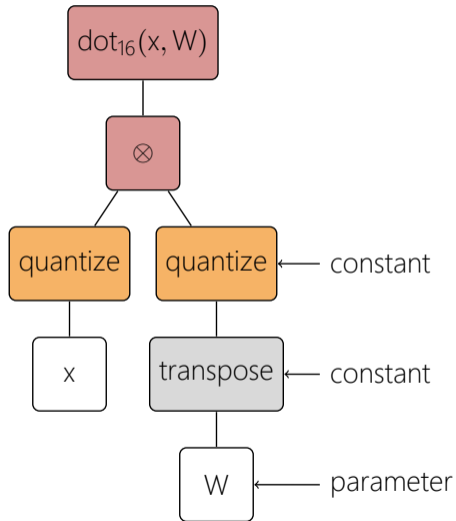
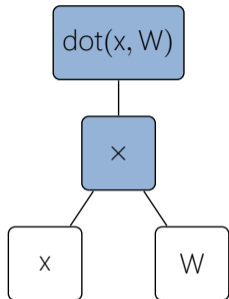


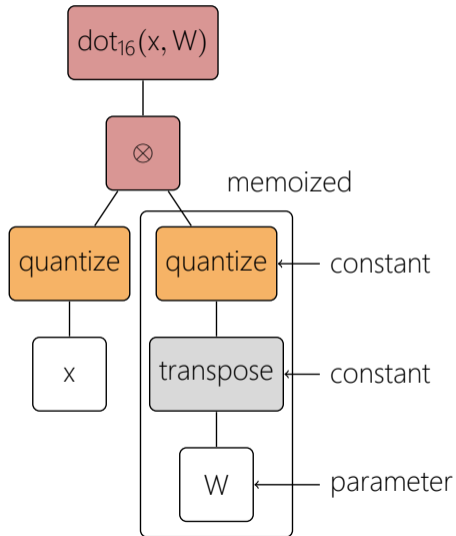
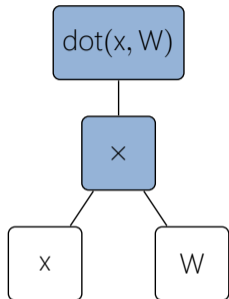
Seconds to translate newstest2014 (batch-size: ca. 384 words – 5 to 25 sentences)

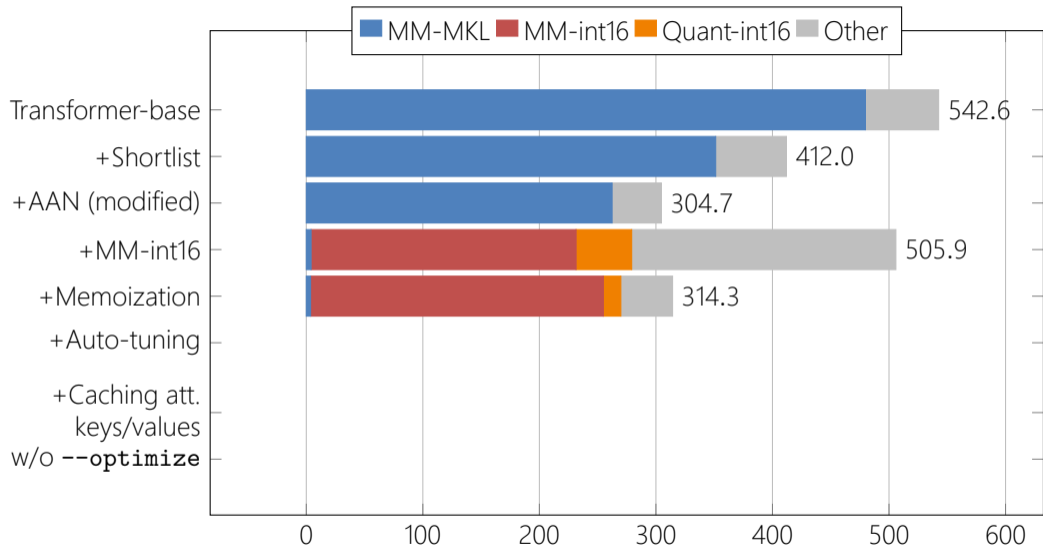




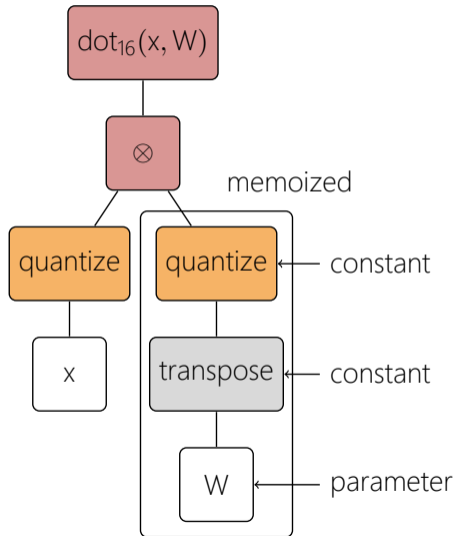
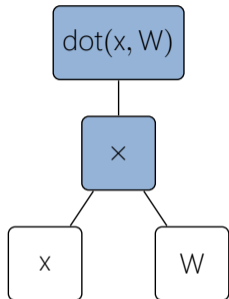


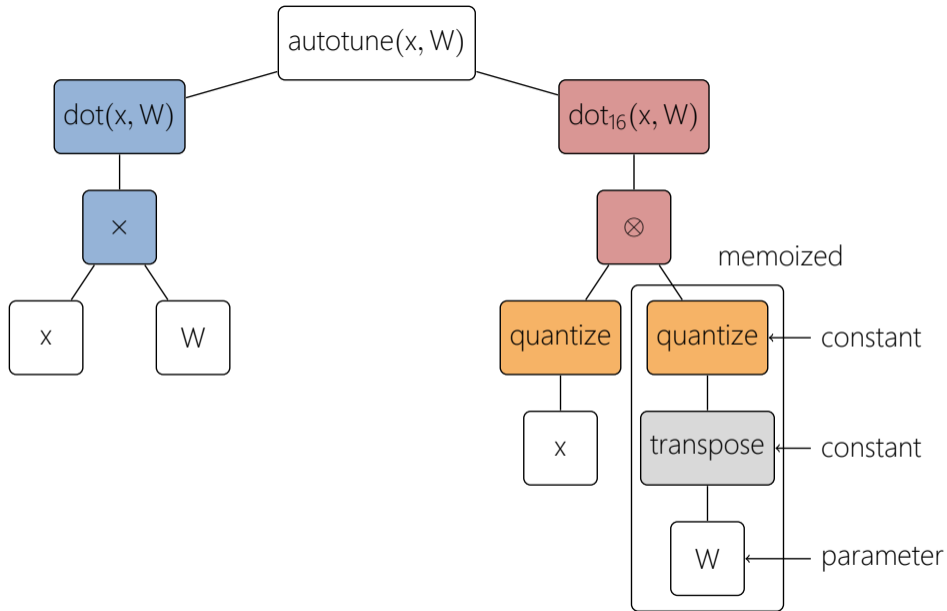






Seconds to translate newstest2014 (batch-size: ca. 384 words – 5 to 25 sentences)

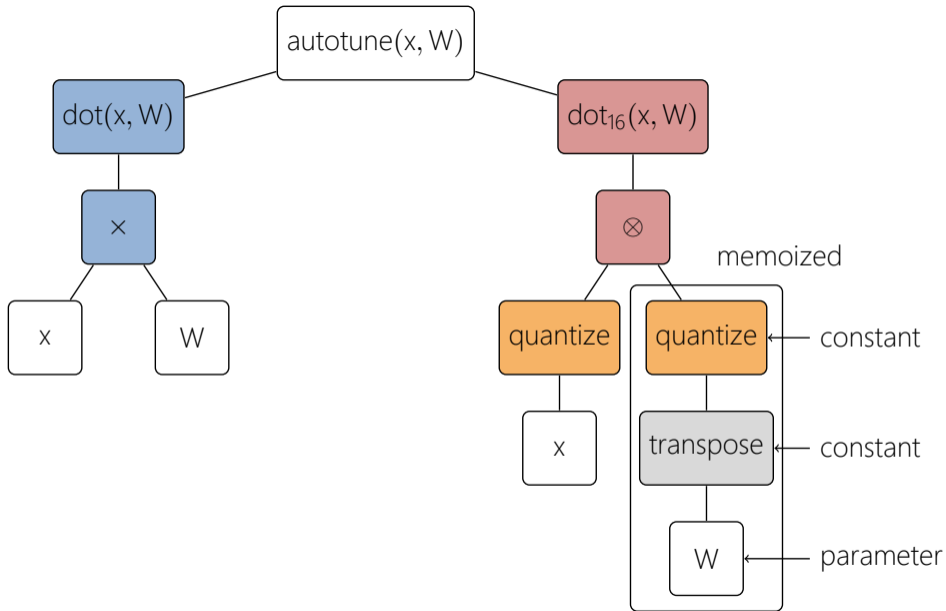


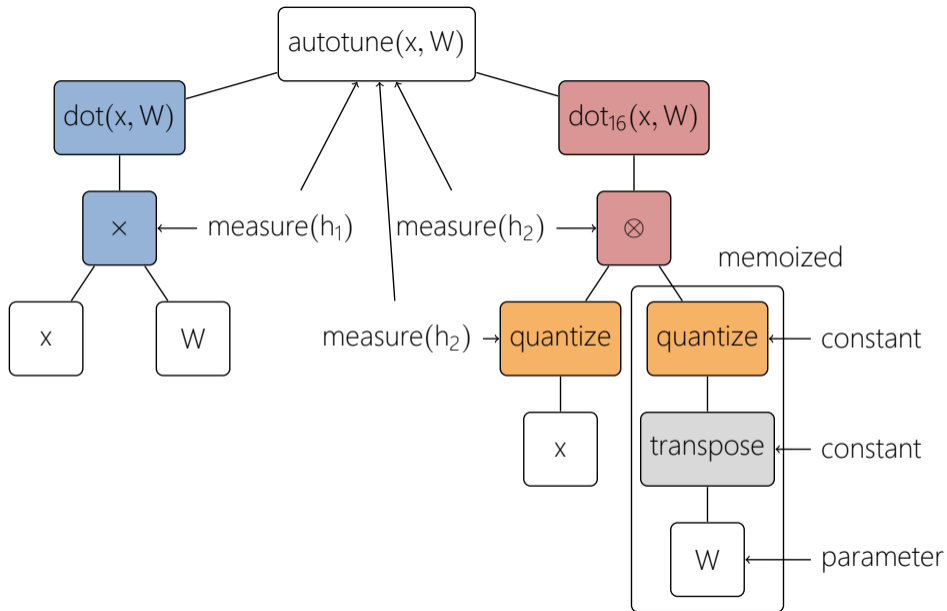


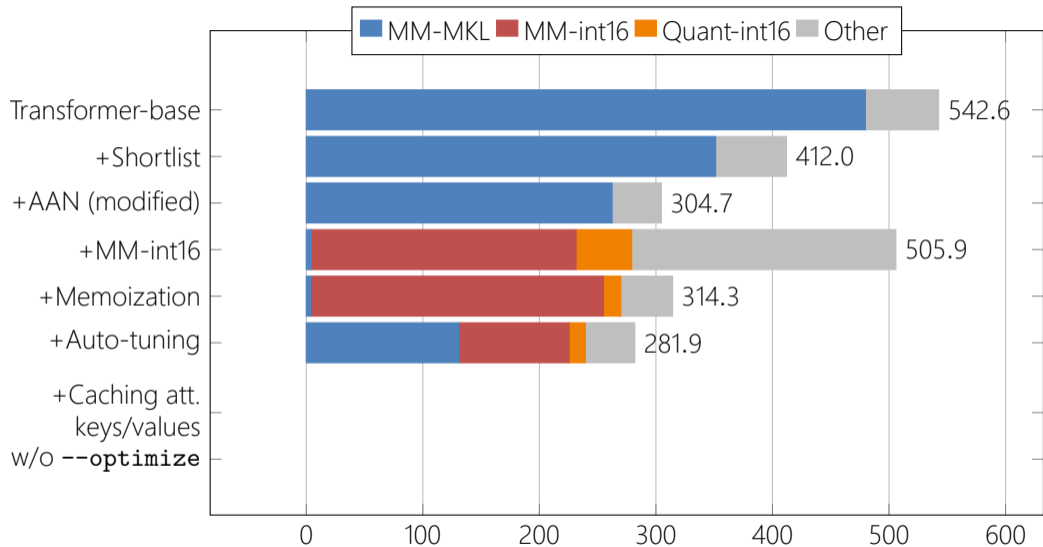
$$\begin{aligned}h_1 &= \text{hash}(\text{dot}(x, W)) \\ &= \text{hash}(\text{dot}) \odot \text{hash}(\text{dims}(x)) \odot \text{hash}(\text{dims}(W)) \\ &= \text{hash}(\text{dot}) \odot \text{hash}(\{11, 512\}) \odot \text{hash}(\{512, 512\})\end{aligned}$$

$$\begin{aligned}h_2 &= \text{hash}(\text{dot}_{16}(x, W)) \\ &= \text{hash}(\text{dot}_{16}) \odot \text{hash}(\text{dims}(x)) \odot \text{hash}(\text{dims}(W)) \\ &= \text{hash}(\text{dot}_{16}) \odot \text{hash}(\{11, 512\}) \odot \text{hash}(\{512, 512\})\end{aligned}$$

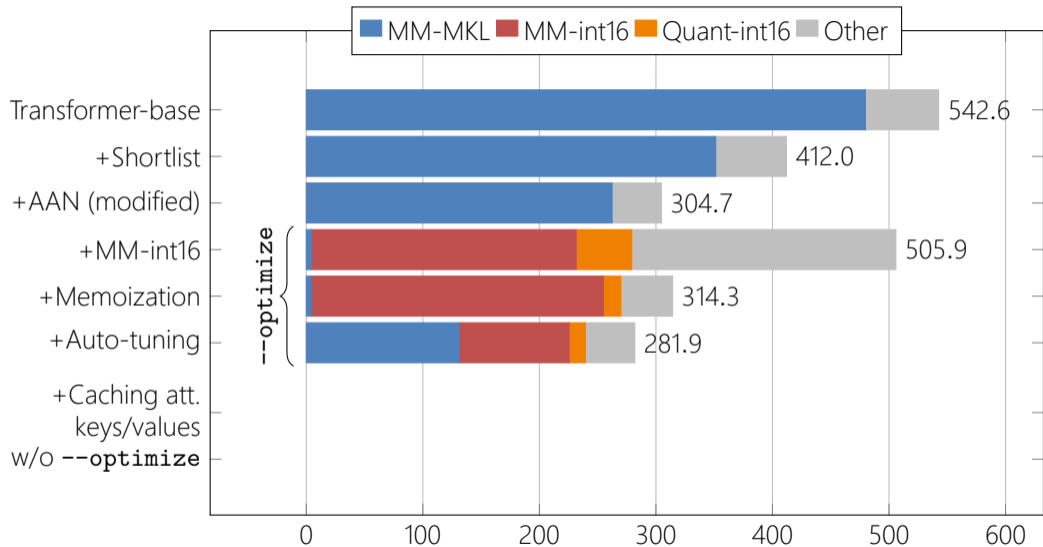
We can decrease granularity via integer-dividing dimensions by a given factor, we choose 5.



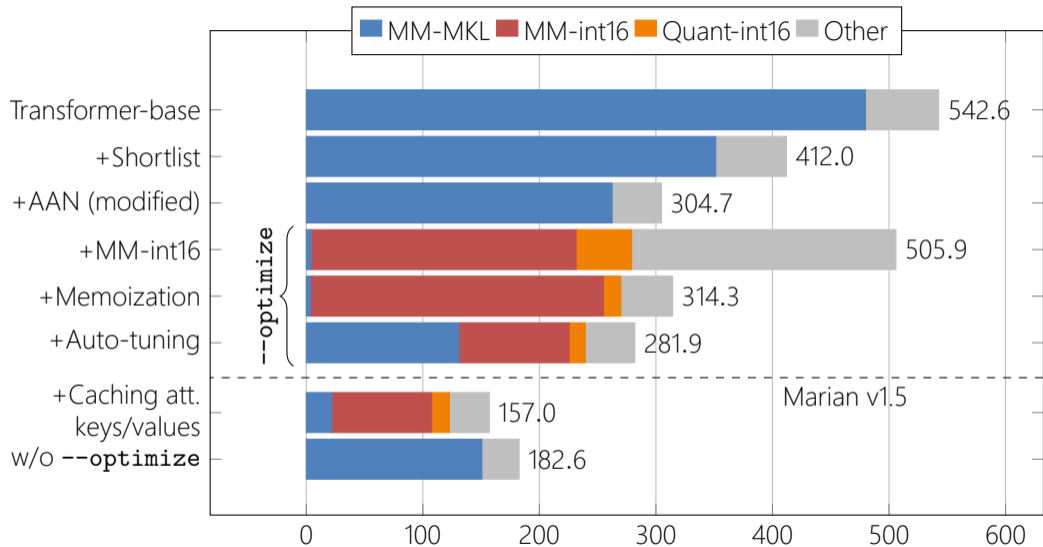




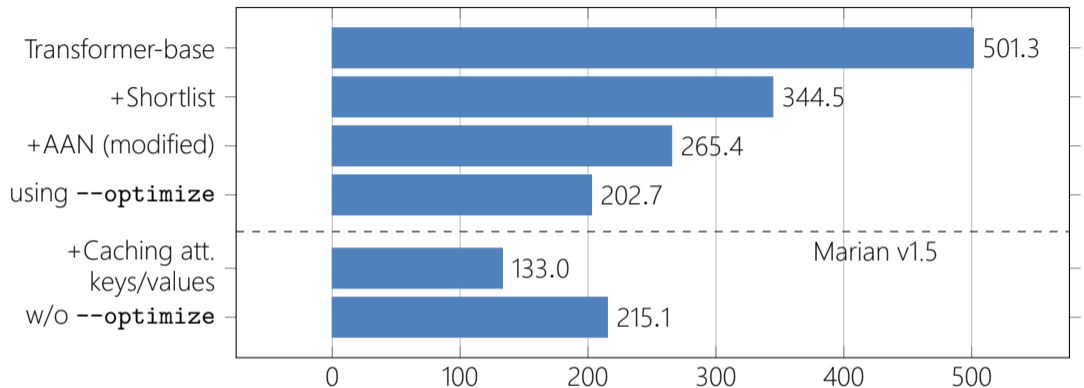
Seconds to translate newstest2014 (batch-size: ca. 384 words – 5 to 25 sentences)



Seconds to translate newstest2014 (batch-size: ca. 384 words – 5 to 25 sentences)



Seconds to translate newstest2014 (batch-size: ca. 384 words – 5 to 25 sentences)



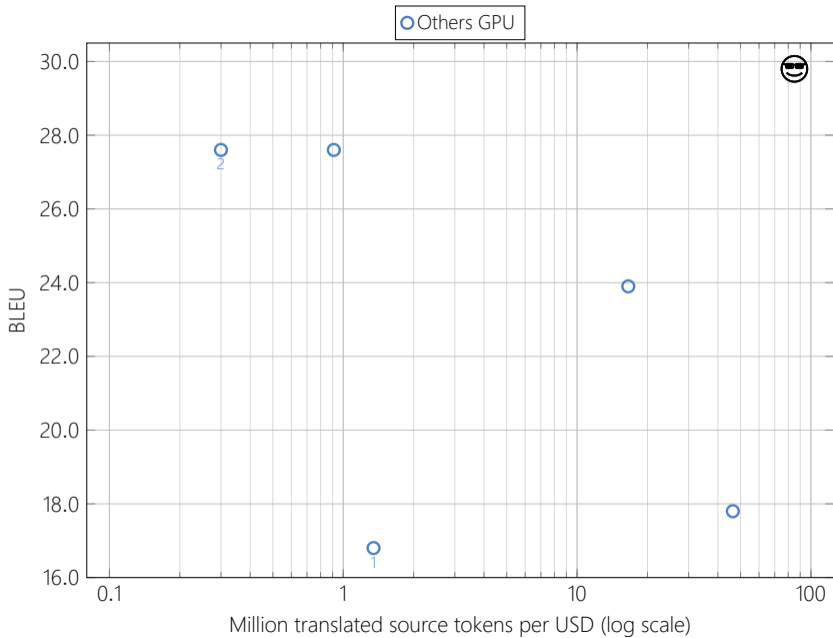
Latency per sentence in milliseconds for newstest2014 (batch-size: 1 sentence)

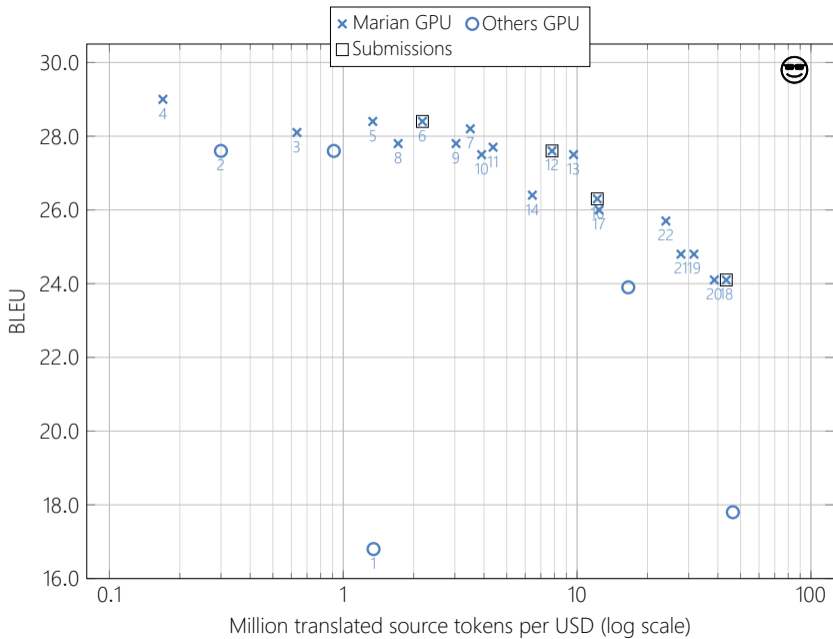
Map GPU and CPU performance into comparable space [w/\$]

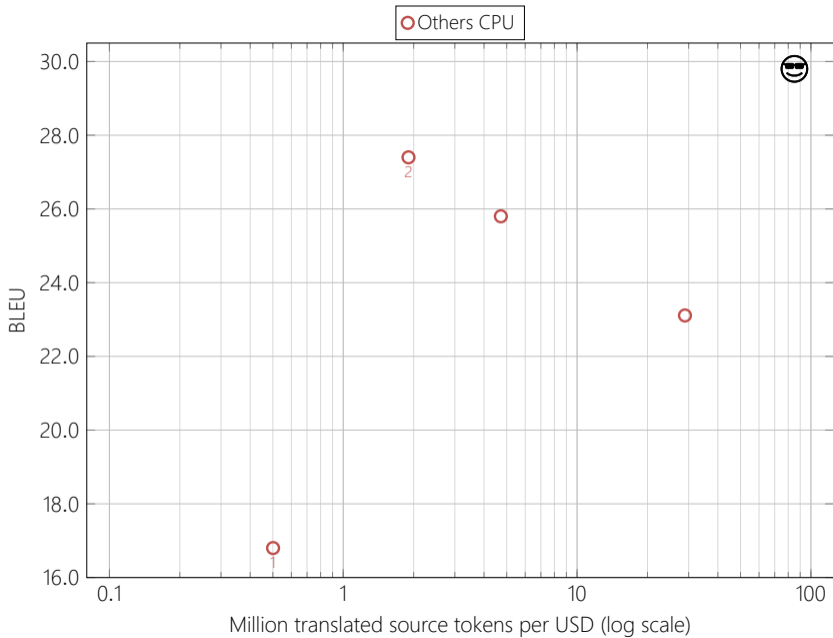
newstest2014.de consists of 62,954 tokens

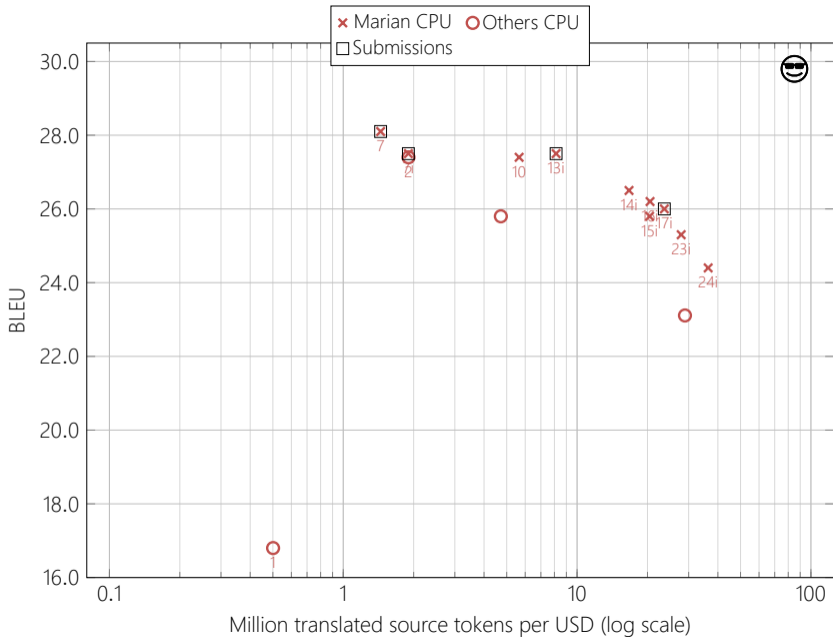
Type	Price [\$/h]
p3.x2large	3.259
m5.large	0.102

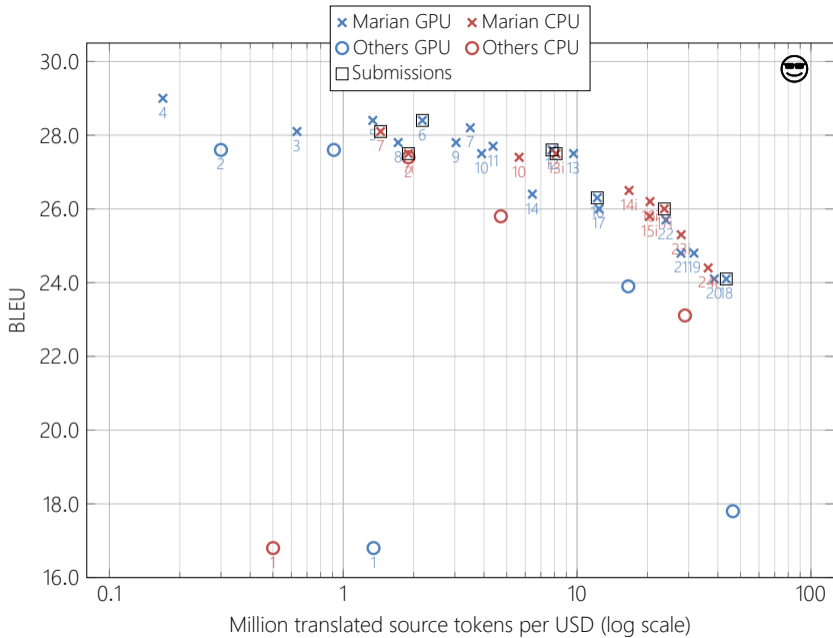
$$[w/\$] = \frac{62,954 [w]}{\text{Translation time [s]}} \cdot \frac{3,600 [s/h]}{\text{Instance price [$/h]}}$$

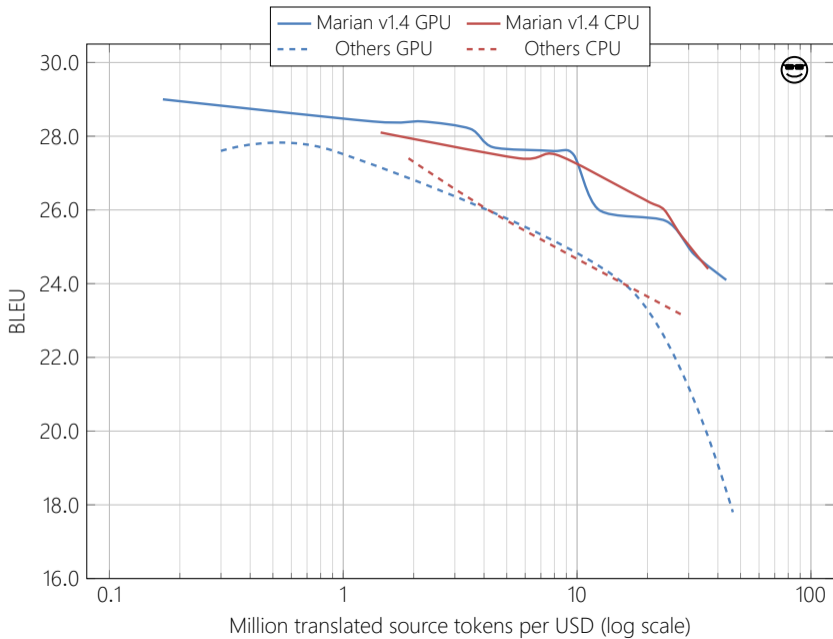


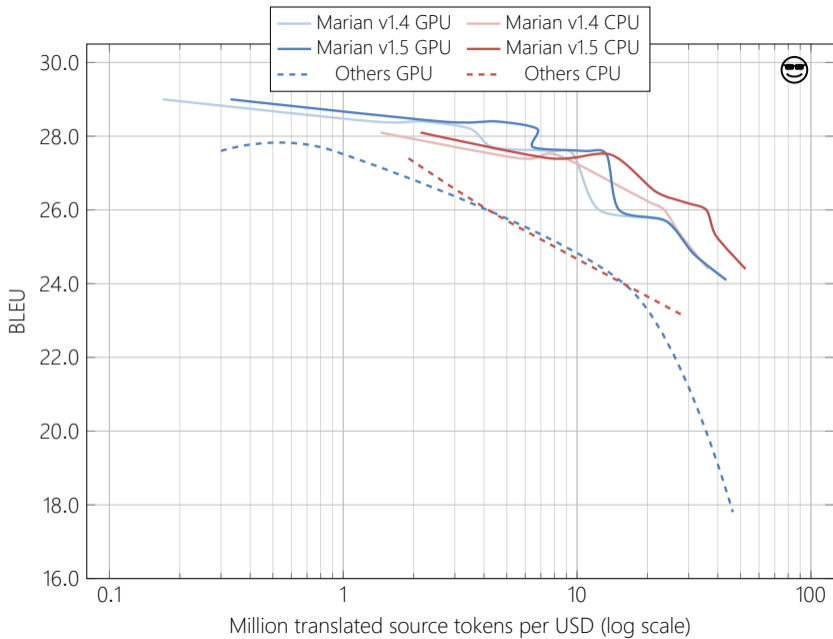












Future work

- More experiments with Teacher-Student scenario;
- More SIMD operations on the CPU;
- All operations in fixed-point arithmetics on the CPU;
- 8-bit matrix product on the CPU;
- Mixed precision (FP16) on the GPU;
- Optimize beam-search for batched translation;
- ...

We are hiring!

Talk to me if you are interested.

Translator