# Exploring Hyper-Parameter Optimization for Neural Machine Translation on GPU Architectures

Robert Lim[1], Kenneth Heafield[2,3], Hieu Hoang[3], Mark Briers[3], and Allen Malony[1]

[1]Computer and Information Science, University of Oregon, Eugene, OR, U.S.A.
[2]School of Informatics, University of Edinburgh, Edinburgh, Scotland, U.K.
[3]The Alan Turing Institute, London, England, U.K.

*Abstract*—**Neural machine translation (NMT) has been accelerated by deep learning neural networks over statistical-based approaches, due to the plethora and programmability of commodity heterogeneous computing architectures such as FPGAs and GPUs and the massive amount of training corpuses generated from news outlets, government agencies and social media. Training a learning classifier for neural networks entails tuning hyper-parameters that would yield the best performance. Unfortunately, the number of parameters for machine translation include discrete categories as well as continuous options, which makes for a combinatorial explosive problem. This research explores optimizing hyper-parameters when training deep learning neural networks for machine translation. Specifically, our work investigates training a language model with Marian NMT. Results compare NMT under various hyper-parameter settings across a variety of modern GPU architecture generations in single node and multi-node settings, revealing insights on which hyper-parameters matter most in terms of performance, such as words processed per second, convergence rates, and translation accuracy, and provides insights on how to best achieve high-performing NMT systems.**

## I. INTRODUCTION

The rapid adoption of neural network (NN) based approaches to machine translation (MT) has been attributed to the massive amounts of datasets, the affordability of high-performing commodity computers, and the accelerated progress in fields such as image recognition, computational systems biology and unmanned vehicles. Research activity in NN-based machine translation has been taking place since the 1990s, but statistical machine translation (SMT) soared along with the successes of machine learning. SMT incorporates a rule-based, data driven approach, and includes language models such as word based (n-grams), phrased-based, syntax-based and hierarchical based approaches. Neural machine translation (NMT), on the other hand, does not require predefined rules, but learns lingusitic rules from statistical models, sequences and occurences from large corpuses. Models trained using NNs produce even higher accuracy than existing SMT approaches, but training time can take anywhere from days to weeks to complete. Suboptimal strategies are often difficult

roblim1@cs.uoregon.edu
kheafiel@inf.ed.ac.uk
hieuhoang@gmail.com
mbriers@turing.ac.uk
malony@cs.uoregon.edu

to find, given the dimensionality and its effect on parameter exploration.

Training a neural network involves the estimation of a huge number of parameters. Ideally, optimization seeks to find the global optima, but in non-convex problems such as neural networks, global optimality is given up and local minima in the parameter space is considered sufficient to obtain models that will generalize beyond training data. Besides obtaining better performance, choosing an appropriate optimization strategy could accelerate the training phase of neural networks and provide higher translation accuracy.

Modeling and training problems are of utmost importance in neural machine translation systems. This work empirically investigates a combination of optimization methods to train a NMT system. The following concerns are addressed: translation performance, training stability, and convergence speed. Specifically, we investigate how well, fast and stable different optimization algorithms are able to find an appropriate local minima, as well as how a combination of these optimizations can solve aspects of problems that arise in model training. Results demonstrate that applying a combination of optimizations leads to faster convergence, translation performance boost and more regularized behavior, compared to selecting an optimizer in isolation.

## II. RELATED WORK

There are many efforts where researchers interpret the characteristics of different optimization techniques. Moreover, other efforts try to show the performance of optimizers in the investigation of loss surface for image classification tasks. Bergstra, et. al. discuss various techniques for hyper-parameter tuning and search strategies, concluding that random search outperforms grid search [1]. Likewise, the authors in [2], [3] take a Bayesian approach toward parameter estimation and optimization. However, these efforts apply their strategies on image recognition tasks.

Britz, et. al. study a massive analysis of NMT hyper-parameters aiming for better optimization being robust to the hyper-parameter variations [4]. Likewise, Bahar et. al. compare various optimization strategies for NMT [5]. In addition, Wu, et. al. [6] utilized the combination of Adam and a simple stochastic gradient descent (SGD) learning algorithm. They run Adam for a fixed number of iterations and switch to SGD to slow down the training phase.

$f$ = (La, croissance, économique, s'est, ralentie, ces, dernières, années, .)

**u**$_i$ — Word Ssample

**p**$_i$ — Word Probability

**z**$_i$ — Recurrent State

**h**$_i$ — Recurrent State

**s**$_i$ — Continuous-space Word Representation

**w**$_i$ — 1-of-K coding

Decoder

Encoder

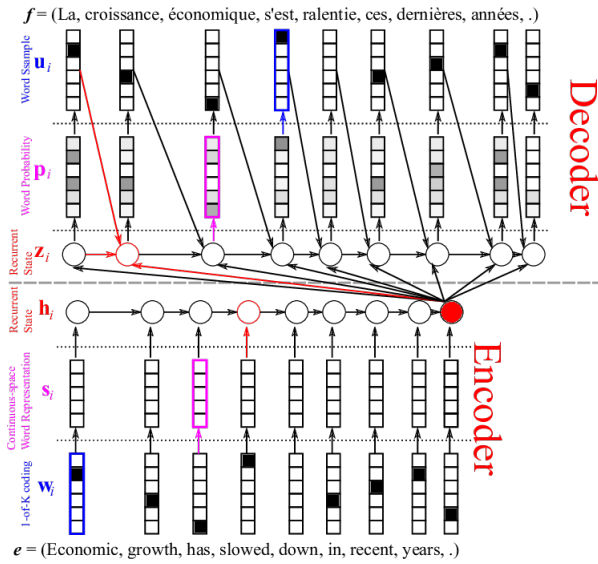$e$ = (Economic, growth, has, slowed, down, in, recent, years, .)

Fig. 1. RNN encoder-decoder, illustrating a sentence translation from English to French. The architecture includes a word embedding space, a 1-of-K coding and a recurrent state on both ends.[2]

To the best of our knowledge, there has not been any work comparing different optimization strategies for NMT. Most of the work in this area focuses on the modeling problem on a vanilla NMT task without exploring the tradeoffs of parameter selection, in terms of performance and stability.

## III. BACKGROUND

Machine translation involves model design and model training. In general, learning algorithms are viewed as a combination of selecting a model criterion (family of functions) and training (parameterization), and a procedure for appropriately optimizing this criterion. The next subsections discuss how sentences are represented with a neural network and the optimization objectives used for training a model for a translation system.

### A. Neural Machine Translation

Given a source $f = f_1^J$ and a target $f = e_1^j$ sequence, NMT models the conditional probability of target words given the source sequence [7]. The NMT training objective function is to minimize the cross-entropy over $S$ training samples $\{\langle f^{(s)}, e^{(s)} \rangle\}_{s=1}^S$ which is defined as follows:

$$J(\theta) = \sum_{s=1}^{S} \sum_{i=1}^{I^{(s)}} \log p(e_i^s | e_{<i}^{(s)}; \theta)$$

Since computing the objective function for the whole training data is expensive, the mini-batch approach randomly selects a small number of samples and averages those samples, resulting in mini-batch gradient calculations.

[2]https://devblogs.nvidia.com/introduction-neural-machine-translation-gpus-part-2/

*1) Recurrent Neural Networks:* A recurrent neural network (RNN) is a neural network with **h** hidden states and an optional $y$ output, and operates on a variable length sequence $X = (x_1, ..., x_T)$. At each time step $t$, the hidden state $\mathbf{h}_{\langle t \rangle}$ of the RNN is updated by

$$\mathbf{h}_{\langle t \rangle} = f(\mathbf{h}_{\langle t-1 \rangle}, x_t), \tag{1}$$

where $f$ is a non-linear activation function, such as a sigmoid or LSTM, and will be discussed in Section III-B2. RNNs can learn a probability distribution over a sequence by being trained to predict the next symbol in a sequence. The output at each timestep $t$ is the conditional distribution $p(x_t | x_{t-1}, ..., x_1)$. For instance, a multinomial distribution (1-of-K coding) can be outputted with the softmax activation function

$$p(x_{t,j} = 1 | x_{t-1}, ..., x_1) = \frac{\exp(\mathbf{w}_j \mathbf{h}_{\langle t \rangle})}{\sum_{j'=1}^{K} \exp(\mathbf{w}_{j'} \mathbf{h}_{\langle t \rangle})} \tag{2}$$

for all possible symbols $j = 1, ..., K$, where $\mathbf{w}_j$ represents the rows of weight matrix $\mathbf{W}$. Thus, the result is a combination of probabilities to compute the probability of sequence $x$ using

$$p(\mathbf{x}) = \prod_{t=1}^{T} p(x_t | x_{t-1}, ..., x_1). \tag{3}$$

The learned distribution is then used to sample a new sequence by iteratively sampling a symbol at each time step. The conditional distribution of the next symbol is defined as

$$P(y_t | y_{t-1}, y_{t-2}, ..., y_1, \mathbf{c}) = g(\mathbf{h}_{\langle t \rangle}, y_{t-1}, \mathbf{c})$$

Note that for activation functions $f$ and $g$, $g$ must produce valid probabilities, where a softmax is typically employed, as defined in Eq. 2.

*2) RNN Encoder-Decoder:* A RNN encoder-decoder (pictured in Fig. 1) encodes a variable-length sequence into a fixed-length vector representation, and decodes a fixed-length vector back into a variable length sequence [8]. The two components of the RNN encoder-decoder are jointly trained to maximize the conditional log-likelihood

$$\arg \max_{\theta} \frac{1}{N} \sum_{n=1}^{N} \log p_{\theta}(y_n | x_n) \tag{4}$$

where $\theta$ represents the set of model parameters, each $\mathbf{x}_n, \mathbf{y}_n$ is a pair of input and output sequences from the training set, and the output of the decoder starting from the input is differentiable. Gradient-based algorithms are used to estimate the model parameters, discussed in Sec III-B. A trained RNN encoder-decoder model can be used to generate a target sequence given an input sequence, and score a given pair of input-output sequence, where the score is simply $p_{\theta}(\mathbf{y}|\mathbf{x})$ from Eqs 3 and 4.

In typical machine translation systems, the goal of the decoder is to find a translation $f$ given source sentence $e$ that maximizes

$$p(\mathbf{f}|\mathbf{e}) \propto p(\mathbf{e}|\mathbf{f})p(\mathbf{f})$$

where $p(\mathbf{e}|\mathbf{f})$ is the translation model, and $p(\mathbf{f})$ represents the language model. In practice, most systems are modeled $\log p(\mathbf{f}|\mathbf{e})$ as a log-linear model with additional features and corresponding weights:

$$\log p(\mathbf{f}|\mathbf{e}) = \sum_{n=1}^{N} w_n f_n(\mathbf{f}, \mathbf{e}) + \log Z(\mathbf{e})$$

where $f_n$ and $w_n$ are the $n^{th}$ feature and weight, and $Z(e)$ is the normalization constant that does not depend on weights. The weights are optimized to maximize the BLEU score on the development set.

SGD

$$\mathbf{g}_t \leftarrow \triangledown_{\theta_t} J(\theta_t)$$
$$\theta_{t+1} \leftarrow \theta_t - \eta \mathbf{g}_t$$

Adagrad

$$\mathbf{g}_t \leftarrow \triangledown_{\theta_t} J(\theta_t)$$
$$\eta_t \leftarrow \eta_{t-1} + \mathbf{g}_t^2$$
$$\theta_{t+1} \leftarrow \theta_t - \frac{\eta}{\sqrt{\eta_t} + \epsilon} \mathbf{g}_t$$

Adam

$$\mathbf{g}_t \leftarrow \triangledown_{\theta_t} J(\theta_t)$$
$$\eta_t \leftarrow \gamma\eta_{t-1} + (1-\gamma)\mathbf{g}_t^2$$
$$\hat{\eta} \leftarrow \frac{\eta_t}{1-\gamma^t}$$
$$\mathbf{m}_t \leftarrow \mu\mathbf{m}_{t-1} + (1-\mu)\mathbf{g}_t$$
$$\hat{\mathbf{m}} \leftarrow \frac{\mathbf{m}_t}{1-\mu^t}$$
$$\theta_{t+1} \leftarrow \theta_t - \frac{\eta}{\sqrt{\hat{\eta}_t} + \epsilon} \hat{\mathbf{m}}_t$$
$$(5)$$

### B. Optimization Objectives

The following subsections describe the tuning of hyperparameters that affect the performance of training a NMT system.

*1) SGD Optimizers:* Stochastic gradient descent (SGD) is commonly used to train neural networks. SGD updates a set of parameters $\theta$, where $\eta$ is the learning rate, or how large the update should be, and $\mathbf{g}_t$ represents the gradient cost function $J(\cdot)$. SGD uses a scheduling-based step size selection, which makes the learning rate $\eta$ an important hyperparameter that requires careful tuning. An adaptive optimizer separately adapts the learning rate for each parameter.

Adaptive moment estimation (Adam) accumulates the decaying average of past squared gradients $\eta_t$. Similar to RmsProp and AdaDelta, Adam also stores the decaying mean of past gradients $\mathbf{m}_t$. The moments, $\hat{m}_t, \hat{n}_t$ are biased corrected terms for instability against zero initialization. Equation 5 displays SGD, AdaGrad and Adam optimizers.

*2) Activation Functions:* To address the vanishing gradients problem associated with learning long-term dependencies in RNNs, LSTMs [9] and GRUs [8] employ a gating mechanism when computing the hidden states. Equation 6 displays the LSTM and GRU activation functions.

For LSTMs (Eq. 6, left), note that the input $i$, forget $f$ and hidden $h$ gates are the same equations except with different parameter matrices, which represent gates that are squashed by the sigmoid into vectors between 0 and 1 values. Multiplying the vectors determine how much of the other vectors to let into the current input state. $g$ is a candidate hidden state that

TABLE I
DATASETS USED IN EXPERIMENTS.

| | RO→EN, EN→RO | DE→EN, EN→DE |
|---|---|---|
| Train | corpus.bpe (2603030) | corpus.bpe (4497879) |
| Valid | newsdev2016.bpe (1999) | newstest2014.bpe (3003) |
| Test | newstest2016.bpe (1999) | newstest2016.bpe (2999) |

is computed based on the current input and previous hidden state. $c_t$ serves as the internal memory, which is a combination of the previous memory $c_{t-1}$ multiplied by the input gate, a balance of two extremes of ignoring either the memory or the new hidden state completely. Lastly, the hidden state $s_t$ calculated as a combination of the internal memory and the output gate.

On the other hand, a GRU (Eq. 6, right) has two gates: a reset gate $r$ and an update gate $u$. The reset gate $r$ determines how to combine the new input with the previous memory, whereas the update gate $u$ defines how much of the previous memory to keep around. If the reset gates were set to 1's and the update gates to 0's, the result would be equivalent to a vanilla RNN.

The differences between the two approaches to compute hidden units are that GRUs have 2 gates, whereas LSTMs have 3 gates. GRUs do not have an internal memory and output gates, compared with LSTM which uses $c$ as its internal memory and $o$ as an output gate. The GRU input and forget gates are coupled by an update gate $z$, and the reset gate $r$ is applied directly to the previous hidden state. Also, there is no 2nd non-linearity in GRUs, compared to LSTMs which uses two hyperbolic tangents.

LSTM

$$i = \sigma(x_t U^i + s_{t-1} W^i)$$
$$f = \sigma(x_t U^f + s_{t-1} W^f)$$
$$o = \sigma(x_t U^o + s_{t-1} W^o)$$
$$g = \tanh(x_t U^g + s_{t-1} W^g)$$
$$c_t = c_{t-1} \circ f + g \circ i$$
$$s_t = \tanh(c_t) \circ o$$

GRU

$$z = \sigma(x_t U^z + s_{t-1} W^z)$$
$$r = \sigma(x_t U^r + s_{t-1} W^r)$$
$$h = \tanh(x_t U^h + (s_{t-1} \circ r) W^h)$$
$$s_t = (1-z) \circ h + z \circ s_{t-1}$$

$$(6)$$

*3) Dropout:* In a fully-connected, feed-forward neural network, dropout randomly retains connections within hidden layers while discarding others [10]. Equation 7 displays a standard hidden update function on the left, whereas a version that decides whether to retain a connection is displayed on the right. $\hat{\mathbf{y}}^{(l)}$ is the thinned output layer, and retaining a network connection is decided by a Bernoulli random variable $r^{(l)}$ with probability $p(\cdot) = 1$.

$$z_i^{(l+1)} = \mathbf{w}_i^{(l+1)}\mathbf{y}^l + b_i^{(l+1)} \qquad r_j^{(l)} \sim \text{Bernoulli}(p) \qquad (7)$$
$$y_i^{(l+1)} = f(z_i^{(l+1)}) \qquad\qquad \hat{\mathbf{y}}^{(l)} = \mathbf{r}^{(l)} * \mathbf{y}^{(l)}$$
$$z_i^{(l+1)} = \mathbf{w}_i^{(l+1)}\hat{\mathbf{y}}^l + b_i^{(l+1)}$$
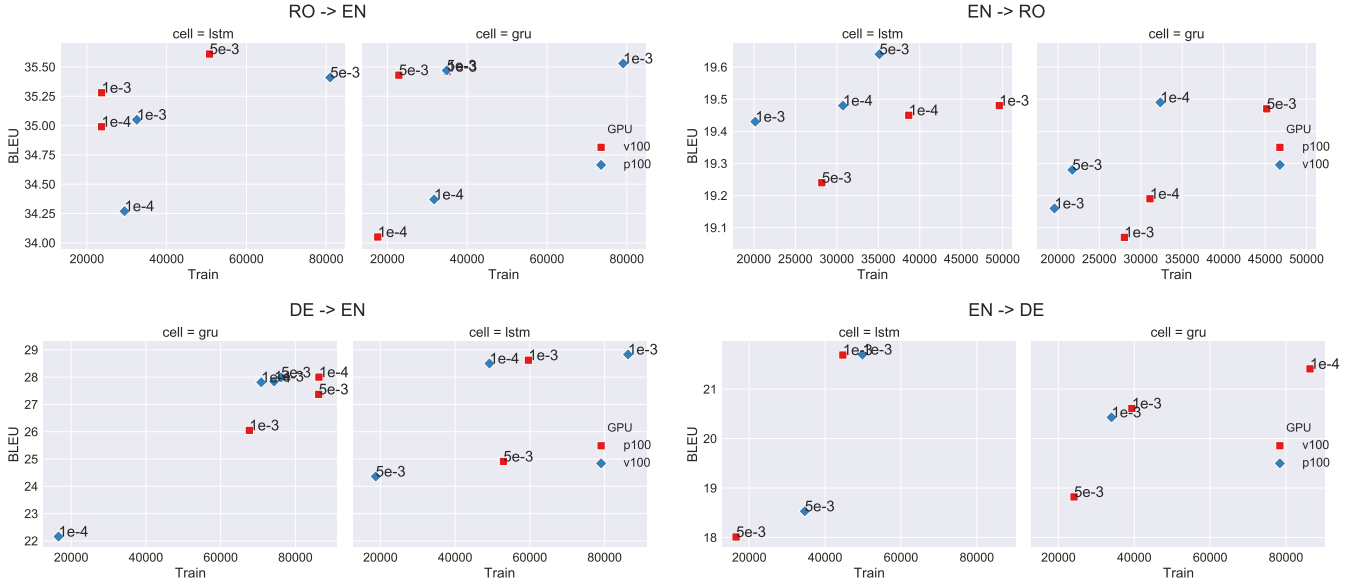$$y_i^{(l+1)} = f(z_i^{(l+1)})$$

Fig. 2. BLEU scores as a function of training time (seconds), comparing GPUs (color), activation units (sub-columns), learning rates and translation directions.
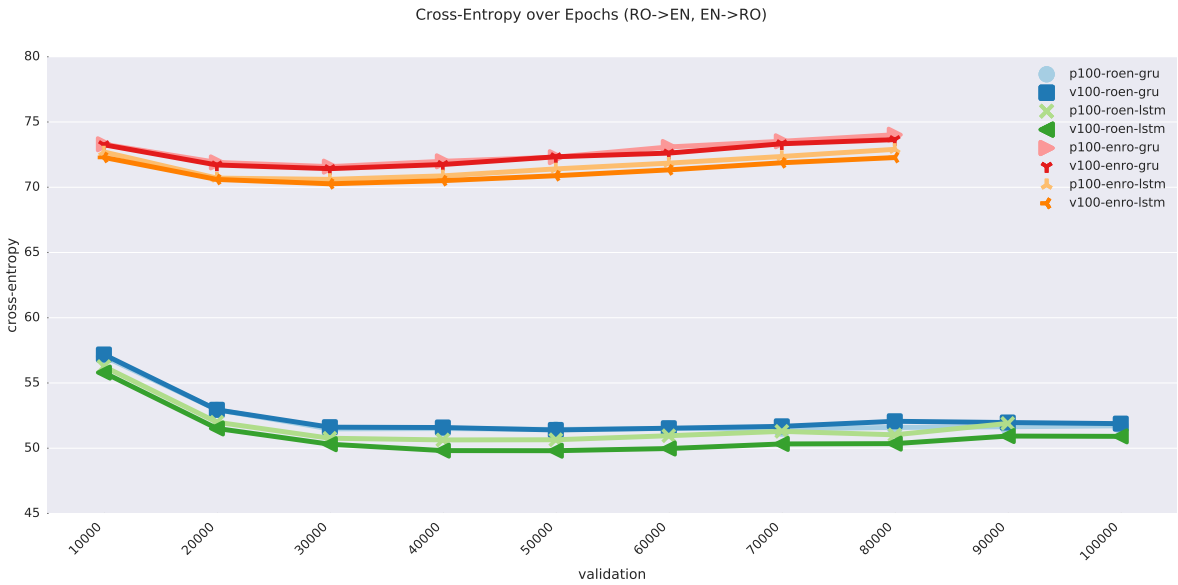


Fig. 3. Cross entropy over the number of epochs for RO → EN and EN → RO, comparing activation functions and GPUs.

## C. Combination of Optimizers

Since the learning trajectory significantly affects the training process, it is required to select and tune the proper types of hyper-parameters to yield good performance. The construction of the RNN cell with activation functions, the optimizer and its learning rate, and the dropout rates all have an affect on how the training progresses, and whether good accuracy can be achieved.

## IV. EXPERIMENTS

The experiments were carried out on the WMT 2016 [11] translation tasks for the Romanian and German languages in four directions: EN → RO, RO → EN, EN → DE, and DE →

EN. The datasets and its characteristics used in the experiments are listed in Table I, with number of sentence examples in parenthesis. Table I shows that for WMT 2016 EN → RO and RO → EN, the training data consisted of 2.6M English and Romanian sentence pairs, whereas for WMT 2016 EN → DE and DE → EN, the training corpus consisted of approximately 4.5M German and English sentence pairs. Validation was performed on 1000 sentences of the newsdev2016 corpus for RO, and on newstest2014 corpus for DE. The newstest2016 corpus consisted of 1999 sentences for RO and 2999 sentences for DE, and was used as the test set. We evaluated and saved the models every 10K iterations and stopped training after 500K iterations.

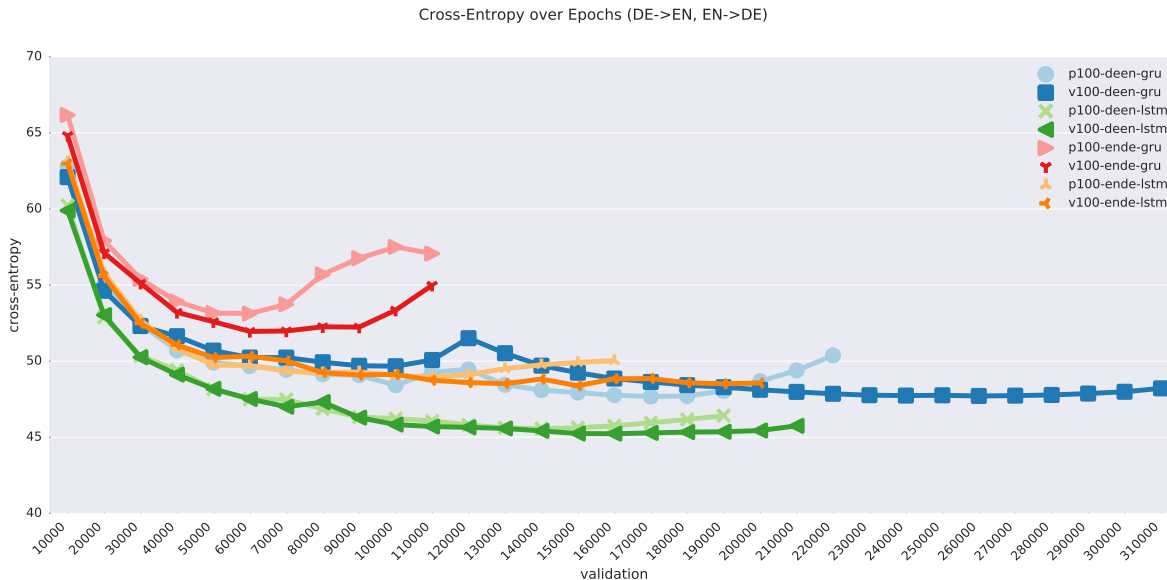All experiments used bilingual data without additional

Fig. 4. Cross-entropy over the number of epochs for DE → EN and EN → DE, comparing activation functions and GPUs.

TABLE II
GRAPHICAL PROCESSORS USED IN THIS EXPERIMENT.

|  | P100 | V100 |
|---|---|---|
| CUDA capability | 6.0 | 7.0 |
| Global memory (MB) | 16276 | 16152 |
| Multiprocessors (MP) | 56 | 80 |
| CUDA cores per MP | 64 | 64 |
| CUDA cores | 3584 | 5120 |
| GPU clock rate (MHz) | 405 | 1380 |
| Memory clock rate (MHz) | 715 | 877 |
| L2 cache size (MB) | 4.194 | 6.291 |
| Constant memory (bytes) | 65536 | 65536 |
| Shared mem blk (bytes) | 49152 | 49152 |
| Registers per block | 65536 | 65536 |
| Warp size | 32 | 32 |
| Max threads per MP | 2048 | 2048 |
| Max threads per block | 1024 | 1024 |
| CPU (Intel) | Ivy Bridge | Haswell |
| Architecture family | Pascal | Volta |

TABLE III
HARDWARE AND EXECUTION ENVIRONMENT INFORMATION.

| Architecture | Haswell | Ivy Bridge |
|---|---|---|
| Model | E5-2698 v3 | Xeon X5650 |
| Clock speed | 2.30 GHz | 2.67 GHz |
| Node count | 4, 14 | 6 |
| GPUs | 4 × V100 | 4 × P100 |
| Memory | 256 GB | 50 GB |
| Linux kernel | 3.10.0-229.14.1 | 2.6.32-642.4.2 |
| Compiler | CUDA v9.0.67 | |
| Flags | {'g', 'lineinfo', 'arch=sm_cc'} | |

monolingual data. The models were trained on Marian [12], an efficient NMT framework written in C++ with multi-node, multi-GPU capabilities. We used the joint byte precision encoding (BPE) approach [13] in both the source and target sets, which converts words to a sequence of subwords. For all four tasks, the number of joint-BPE operations were 20K. All words were projected on a 512-dimensional embedding

space, with vocabulary dimensions of $66000 \times 50000$. The mini-batch size was determined automatically based on the sentence length that was able to fit in GPU global memory, set at 13000 MB for each GPU.

Decoding was performed using beam search with a beam size of 12. The translation portion consisted of recasing and detokenizing the translated BPE chunks. The trained models compared different hyper-parameter strategies, including the type of optimizer, the activation function, and the amount of dropout applied, as discussed in Section III-C. The number of parameters were initialized with the same random seed. The systems were evaluated using the case-sensitive BLEU score computed by Moses SMT [14].

We compared models trained on two different types of GPUs (P100 Pascal, V100 Volta), listed on Table II. The corresponding CPUs are listed on Table III. Each ran with four GPUs. The dataset was partitioned across 4 GPUs, and a copy of the model was executed on each GPU.

## V. ANALYSIS

This section analyzes the results of the evaluated NMT systems in terms of translation quality, training stability and convergence speed.

### A. Translation Qualtiy

Table IV shows BLEU scores calculated for four translation directions for the validation sets (top) and the test sets (bottom), comparing learning rates, activation functions and GPUs. Note that entries with n/a means that no results were available, whereas entries with dnf indicates training time that did not complete within 24 hours. For the validation sets, LSTM was able to achieve higher accuracy rates, whereas in the test set GRUs and LSTMs were about the same. Also, note that the best performing learning rates were usually at a lower value (e.g. 1e-3). The type of hidden unit mechanism

TABLE IV
BLEU SCORES FOR VALIDATION (TOP) AND TEST (BOTTOM) DATASETS.

| cell | learn-rt | ro→en | | en→ro | | de→en | | en→de | |
|------|----------|-------|------|-------|------|-------|------|-------|------|
| | | P100 | V100 | P100 | V100 | P100 | V100 | P100 | V100 |
| GRU | 1e-3 | 35.53 | 35.43 | 19.19 | 19.28 | 28.00 | 27.84 | 20.43 | 20.61 |
| | 5e-3 | 34.37 | 34.05 | 19.07 | 19.16 | 26.05 | 22.16 | n/a | 19.01 |
| | 1e-4 | 35.47 | 35.46 | 19.45 | 19.49 | 27.37 | 27.81 | dnf | 21.41 |
| LSTM | 1e-3 | 34.27 | 35.61 | 19.29 | 19.64 | 28.62 | 28.83 | 21.70 | 21.69 |
| | 5e-3 | 35.05 | 34.99 | 19.48 | 19.43 | n/a | 24.36 | 18.53 | 18.01 |
| | 1e-4 | 35.41 | 35.28 | 19.43 | 19.48 | n/a | 28.50 | dnf | dnf |
| GRU | 1e-3 | 34.22 | 34.17 | 19.42 | 19.43 | 33.03 | 32.55 | 26.55 | 26.85 |
| | 5e-3 | 33.13 | 32.74 | 19.31 | 18.97 | 31.04 | 26.76 | n/a | 26.02 |
| | 1e-4 | 33.67 | 34.44 | 18.98 | 19.69 | 33.15 | 33.12 | dnf | 28.43 |
| LSTM | 1e-3 | 33.10 | 33.95 | 19.56 | 19.08 | 33.10 | 33.89 | 28.79 | 28.84 |
| | 5e-3 | 33.10 | 33.52 | 19.13 | 19.51 | n/a | 29.16 | 24.12 | 24.12 |
| | 1e-4 | 33.29 | 32.92 | 19.14 | 19.23 | n/a | 33.44 | dnf | dnf |

TABLE V
DROPOUT RATES, BLEU SCORES AND TOTAL TRAINING TIME FOR TEST SET, COMPARING SYSTEMS.

| cell | dropout | ro→en | | | | de→en | | | |
|------|---------|-------|------|-------|------|-------|------|-------|------|
| | | P100 | $t$ | V100 | $t$ | P100 | $t$ | V100 | $t$ |
| GRU | 0.0 | 34.47 | 6:29 | 34.47 | 4:43 | 32.29 | 9:48 | 31.61 | 6:15 |
| | 0.2 | 35.53 | 8:48 | 35.43 | 6:21 | 33.03 | 18:47 | 32.55 | 19:40 |
| | 0.3 | 35.36 | 12.21 | 35.15 | 7:28 | 31.36 | 10:14 | 31.50 | 9:33 |
| | 0.5 | 34.50 | 12:20 | 34.67 | 17:18 | 29.64 | 11:09 | 30.21 | 11:09 |
| LSTM | 0.0 | 34.84 | 6:29 | 34.65 | 4:46 | 32.84 | 12:17 | 32.88 | 7:37 |
| | 0.2 | 34.27 | 8:10 | 35.61 | 6:34 | 33.10 | 16:33 | 33.89 | 13:39 |
| | 0.3 | 35.67 | 9:56 | 35.37 | 11:29 | 33.45 | 20.02 | 33.51 | 15:51 |
| | 0.5 | 34.50 | 15:13 | 34.33 | 12:45 | 32.67 | 20:02 | 32.20 | 13:03 |

(e.g LSTM vs GRU) and the learning rate can affect the overall accuracy achieved, as demonstrated by Table IV.

Table V displays various dropout rates applied for two translation directions RO → EN and DE → EN, comparing hidden units, GPUs and overall training time. The learning rate was evaluated at 0.001, the rate that achieved the highest BLEU score as evident in Table IV. Generally speaking, increasing the dropout rates also increased training time. This may be the result of losing network connections when applying the dropout mechanism, but at the added benefit of avoiding overfitting. This is evident in Table V, where applying some form of dropout will result in a trained model achieving higher accuracies. The best performance can be seen when the dropout rate was set at 0.2 to 0.3. This confirms that some form of skip connection mechanism is necessary to prevent the overfitting of models under training.

Figure 2 shows BLEU score results as a function of training time, comparing GPUs, activation units, learning rates and translation directions. Note that in most cases a learning rate of 0.001 achieves the higher accuracy in most cases, at the cost of higher training time. Also, note the correlation between longer training time and higher BLEU scores in most cases. In some cases, the models were able to converge at a faster rate (e.g. Fig. 2 upper left, RO→EN, GRU with learning rate of 0.005 vs 0.001).

### B. Training Stability

Figure 3 shows the cross-entropy scores for the RO → EN and EN → RO translation tasks, comparing different activation functions (GRU vs. LSTM), with learning rates at 0.001. Note the training stability patterns that emerge from this plot, which is highly correlated with the translation direction. The activation function (GRU vs LSTM) during validation also performed similarly across GPUs and was also highly correlated with the translation direction. Cross-entropy scores for the EN → RO translation direction were more or less the same. However, for RO → EN, a LSTM that executed on a P100 converged the earliest by one iteration.

Figure 4 shows the same comparison of cross-entropy scores over epochs for DE → EN and EN → DE translation tasks. Note that the behavior for this translation task was wildly different for all systems. Not only did it take more epochs to converge compared to Fig 3, but also how well the system progressed also varies, as evident in the cross-entropy scores during validation. When comparing hidden units, LSTMs outperformed GRUs in all cases. When comparing GPUs, the V100 performed better than the P100 in terms of cross-entropy, but took longer to converge in some cases (e.g. v100-deen-lstm, v100-ende-lstm). Also, note that the behavior of the translation task EN → DE for a GRU hidden unit never stabilized, as evident in both the high cross-entropy scores and the peaks toward the end. The LSTM was able to achieve a better cross-entropy score overall, with nearly a 8 point difference for DE → EN, compared with the GRU.

### C. Convergence Speed

Figure 5 shows the average words-per-second for the RO → EN translation task, comparing systems. The average words-per-second executed remained consistent across epochs. The system that was able to achieve the most words-per-second was v100-roen-gru-0.001, whereas the one that achieved the
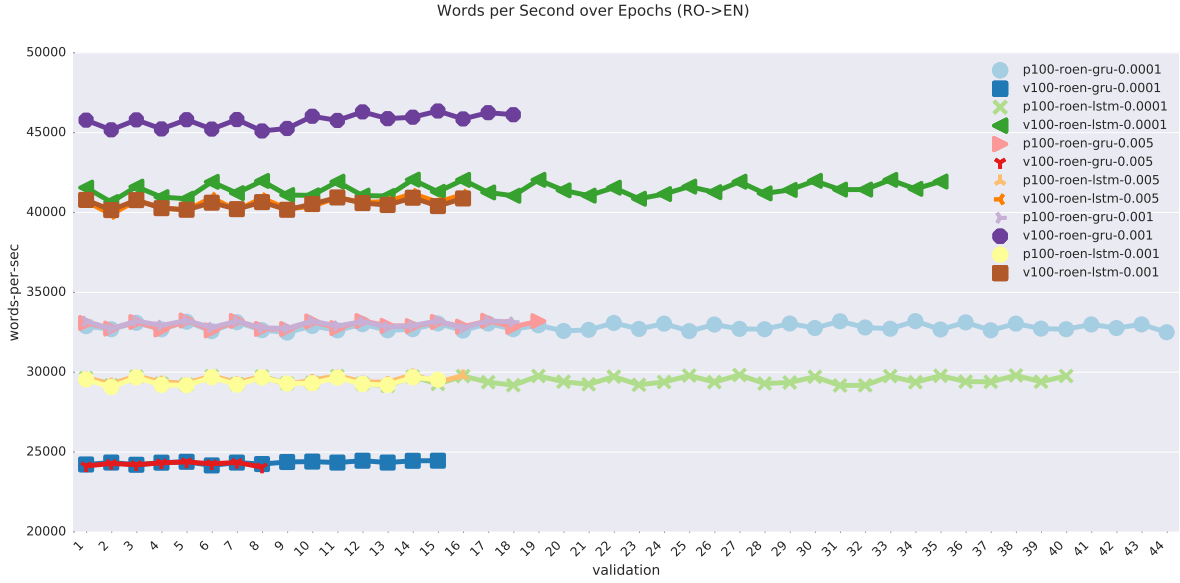
Fig. 5. Average words-per-second for the RO → EN translation task, comparing systems.

least words-per-second was the v100-roen-gru-0.005. Surprisingly, the best and worst performer was the v100-roen-gru, depending on its learning rate, with the sweet spot at 0.001. This confirms 0.001 as the best learn rate that can execute a decent number of words-per-second and achieve a fairly high accuracy, as evident in previous studies, across all systems.

Table VI also displays words-per-second and validation, comparing activation units, learning rates and GPUs. When fixing learning rate, the V100 was able to execute more words-per-second than the P100, and was able to converge at an earlier iteration. When comparing hidden units, GRUs were able to execute higher words per second on a GPU and converge at a reasonable rate (at 18000 iterations) for most learning rates, except for 5e-3. When looking at LSTMs, words-per-second executed on a V100 was similar, although at a higher learning rate it was able to converge at 42000 iterations, but at the cost of longer training time and slower convergence (35000 iterations).

Table VII shows the corresponding total training time for the four translation directions, comparing GPUs, activation units, and learning rates. The dropout rate was set at 0.2, which was the best performer in most cases (Tab V). Table VII shows that the training time increased as the learning rates were decreased. In general, Romanian took a fraction of the time to complete training (usually under 10 hours), whereas German took 18-22 hours to complete training.

## VI. Discussion

The variation in the results, in terms of language translation, hyper-parameters, words-per-second executed and BLEU scores, in addition to the hardware the training was executed on demonstrates the complexity in learning the grammatical structure between the two languages. In particular, the learning rate set for training, the hidden unit selected for the activation function, the optimization criterion and the amount of dropout

applied to the hidden connections all have a drastic effect on overall accuracy and training time. Specifically, we found that a lower learning rate achieved the best performance in terms of convergence speed and BLEU score. Also, we found that the V100 is able to execute more words-per-second than the P100 in all cases. When looking at accuracy as a whole, LSTM hidden units outperformed GRUs in all cases. Lastly, the amount of dropout applied on a network in all cases prevented the model from overfitting and achieve a higher accuracy.

The multidimensionality of hyper-parameter optimization poses a challenge in selecting the architecture design for training NN models, as illustrated by the varying degrees of behavior across systems and its performance outcome. This work investigated how the varying design decisions can affect training outcome and provides neural network designers how to best look at which parameters affect performance, whether accuracy, words processed per second, and convergence expectation. Coupled with massive datasets for parallel text corpuses and commodity heterogenous GPU architectures, the models trained were able to achieve WMT grade accuracy with the proper selection of hyper-parameter tuning.

## VII. Conclusion

We analyzed the performance of various hyper-parameters for training a NMT, including the optimization strategy, the learning rate, the activation cell, and the GPU across various systems for the WMT 2016 translation task in four translation directions. Results demonstrate that a proper learning rate and a minimal amount of dropout is able to prevent overfitting as well as achieve high training accuracy.

Future work includes developing optimization methods to evaluate how to best select hyper-parameters. By statically analyzing the computational graph that represents a NN in

TABLE VI
WORDS-PER-SECOND (AVERAGE) AND NUMBER OF EPOCHS, COMPARING ACTIVATION UNITS, LEARNING RATES AND GPUS.

| cell | learn-rt | words-per-sec | | validation | | words-per-sec | | validation | |
|---|---|---|---|---|---|---|---|---|---|
| | | P100 | V100 | P100 | V100 | P100 | V100 | P100 | V100 |
| | | ro→en | | | | en→ro | | | |
| GRU | 1e-3 | 33009.23 | 45762.54 | 18000 | 18000 | 29969.14 | 42746.15 | 15000 | 15000 |
| | 5e-3 | 32965.23 | 24253.14 | 19000 | 8000 | 30223.89 | 23144.62 | 17000 | 10000 |
| | 1e-4 | 32828.61 | 24341.96 | 44000 | 16000 | 29959.34 | 23277.51 | 25000 | 14000 |
| LSTM | 1e-3 | 29412.87 | 40534.06 | 15000 | 16000 | 27282.54 | 38131.13 | 14000 | 14000 |
| | 5e-3 | 29536.65 | 40598.24 | 16000 | 16000 | 27245.42 | 37384.46 | 19000 | 21000 |
| | 1e-4 | 29478.51 | 41441.37 | 40000 | 35000 | 27002.60 | 38118.79 | 25000 | 25000 |
| | | de→en | | | | en→de | | | |
| GRU | 1e-3 | 28279.53 | 38026.87 | 20000 | 28000 | 28367.91 | 39995.48 | 10000 | 10000 |
| | 5e-3 | 28215.40 | 19819.59 | 25000 | 4000 | n/a | 39944.10 | n/a | 16000 |
| | 1e-4 | 28367.54 | 33218.70 | 26000 | 32000 | dnf | 39993.89 | dnf | 36000 |
| LSTM | 1e-3 | 24995.64 | 33507.31 | 16000 | 17000 | 25245.67 | 35122.54 | 13000 | 17000 |
| | 5e-3 | 25210.15 | 33740.92 | 14000 | 7000 | 25049.21 | 33649.20 | 9000 | 6000 |
| | 1e-4 | dnf | 34529.58 | dnf | 31000 | dnf | dnf | dnf | dnf |

TABLE VII
TOTAL TRAINING TIME FOR FOUR TRANSLATION DIRECTIONS, COMPARING SYSTEMS.

| cell | learn-rt | ro→en | | en→ro | | de→en | | en→de | |
|---|---|---|---|---|---|---|---|---|---|
| | | P100 | V100 | P100 | V100 | P100 | V100 | P100 | V100 |
| GRU | 1e-3 | 8:48 | 6:21 | 7:47 | 5:26 | 18:47 | 19:40 | 9:26 | 6:41 |
| | 5e-3 | 9:41 | 4:52 | 8:38 | 6:02 | 23:57 | 4:36 | n/a | 10:56 |
| | 1e-4 | 21:58 | 9:43 | 12:33 | 8:59 | 23:50 | 21:09 | dnf | 23:58 |
| LSTM | 1e-3 | 8:10 | 6:34 | 7:49 | 5:36 | 16:33 | 13:39 | 13:50 | 12:24 |
| | 5e-3 | 9:02 | 6:34 | 10:44 | 8:32 | n/a | 5:12 | 9:37 | 4:35 |
| | 1e-4 | 22:29 | 14:05 | 13:46 | 9:45 | n/a | 23:57 | dnf | dnf |

terms of instruction operations executed and resource allocation constraints, one could derive execution performance for a given dataset without running experiments.

## REFERENCES

[1] J. S. Bergstra, R. Bardenet, Y. Bengio, and B. Kégl, "Algorithms for hyper-parameter optimization," in *Advances in neural information processing systems*, 2011, pp. 2546–2554.

[2] B. Shahriari, K. Swersky, Z. Wang, R. P. Adams, and N. De Freitas, "Taking the human out of the loop: A review of bayesian optimization," *Proceedings of the IEEE*, vol. 104, no. 1, pp. 148–175, 2016.

[3] J. Snoek, H. Larochelle, and R. P. Adams, "Practical bayesian optimization of machine learning algorithms," in *Advances in neural information processing systems*, 2012, pp. 2951–2959.

[4] D. Britz, A. Goldie, T. Luong, and Q. Le, "Massive exploration of neural machine translation architectures," *arXiv preprint arXiv:1703.03906*, 2017.

[5] P. Bahar, T. Alkhouli, J.-T. Peter, C. J.-S. Brix, and H. Ney, "Empirical investigation of optimization algorithms in neural machine translation," *The Prague Bulletin of Mathematical Linguistics*, vol. 108, no. 1, pp. 13–25, 2017.

[6] Y. Wu, M. Schuster, Z. Chen, Q. V. Le, M. Norouzi, W. Macherey, M. Krikun, Y. Cao, Q. Gao, K. Macherey *et al.*, "Google's neural machine translation system: Bridging the gap between human and machine translation," *arXiv preprint arXiv:1609.08144*, 2016.

[7] D. Bahdanau, K. Cho, and Y. Bengio, "Neural machine translation by jointly learning to align and translate," *arXiv preprint arXiv:1409.0473*, 2014.

[8] K. Cho, B. Van Merriënboer, C. Gulcehre, D. Bahdanau, F. Bougares, H. Schwenk, and Y. Bengio, "Learning phrase representations using rnn encoder-decoder for statistical machine translation," *arXiv preprint arXiv:1406.1078*, 2014.

[9] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural computation*, vol. 9, no. 8, pp. 1735–1780, 1997.

[10] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov, "Dropout: A simple way to prevent neural networks from overfitting," *The Journal of Machine Learning Research*, vol. 15, no. 1, pp. 1929–1958, 2014.

[11] M. Junczys-Dowmunt and R. Grundkiewicz, "Log-linear combinations of monolingual and bilingual neural machine translation models for automatic post-editing," in *ACL*. Berlin, Germany: Association for Computational Linguistics, 2016, pp. 751–758. [Online]. Available: http://www.aclweb.org/anthology/W16-2378

[12] M. Junczys-Dowmunt, T. Dwojak, and H. Hoang, "Is neural machine translation ready for deployment? a case study on 30 translation directions," in *IWSLT*, 2016.

[13] R. Sennrich, B. Haddow, and A. Birch, "Neural machine translation of rare words with subword units," *arXiv preprint arXiv:1508.07909*, 2015.

[14] P. Koehn, H. Hoang, A. Birch, C. Callison-Burch, M. Federico, N. Bertoldi, B. Cowan, W. Shen, C. Moran, R. Zens *et al.*, "Moses: Open source toolkit for statistical machine translation," in *ACL*. Association for Computational Linguistics, 2007, pp. 177–180.