

A Primer in NMTology: What we have understood about NMT

Anonymous EMNLP submission

Abstract

Neural Machine Translation (NMT) has been through great revolutions in recent years. Accompanied with improvements in translation quality are works that attempted to understand the working mechanism of various aspects of the NMT framework. In our paper, we survey those efforts on unveiling the *black box* of the standard NMT framework. To begin with, we briefly introduce the three critical components of the holistic NMT framework; next, we deliver a clear *component-centric* categorization and clean summary of these specific works guided by *frequently-asked* questions (FAQs) that aim at making up *lack* of understanding; finally, we discuss several limitations, future directions and inspirations. We believe this paper could facilitate the community to weave a holistic and clear picture of our current understandings of the standard NMT framework and shed light on its future improvements and developments. Please check this website <https://nmtology.github.io/> for a visual guidance of the FAQs.

1 Introduction

Machine Translation is an extremely challenging task. Statistical Machine Translation (SMT), which models translation in a pipelined manner, was historically one of the popular approaches (Koehn et al., 2003; Chiang, 2005). In the SMT pipeline, each module plays a clear role and is parameterized by a relatively simple model, leading to easy interpretability. Recently, Neural Machine Translation (NMT) framework establishes new state-of-the-art performances (Barrault et al., 2019, 2020). The strengths of NMT come from its strong modeling power with *complex* deep encoder-decoder architecture and *holistic* end-to-end training, which lead to poor interpretability. Consequently, the poor interpretability prevents us from elegantly debugging the model, trusting its outputs, and particularly further improving performance (Ding et al., 2017).

This paper conducts a thorough survey on understanding components of the NMT framework, covering a hundred papers published in recent years. Our survey is *component-centric*, that is, we organize related papers in terms of every NMT component and highlight important questions *frequently-asked w.r.t.* that component. We want our readers to treat this paper as instructional FAQs about understanding the black-box of the NMT framework, so they can quickly zoom into certain question and find the corresponding papers to complement their lack of understanding. §2 briefly introduces every components of the NMT framework, while §3, §4, §5, §6 summarize works in terms of model architecture, training, inference and behavior. In §7, we discuss limitations and future directions.

Related surveys Lertvittayakumjorn and Toni (2021); Danilevsky et al. (2020); Luo et al. (2021) are more general surveys on principles for explainable NLP, as they mainly discuss general desiderata and possible explanation paradigms or frameworks. Sajjad et al. (2021) survey specific neuron-level interpretation methods for NLP models, while Belinkov and Glass (2019) and Sun et al. (2021) focus on surveying a broad range of general techniques and methods for interpreting NLP models on various tasks. The closest work in organization to us might be Rogers et al. (2020). They deliver a thorough survey on research questions, directions, and solutions around large pretrained models. In our paper, we organize research works guided by research questions that are more related to the interpretation and understanding side, so that researchers can gain in-depth insights in various components and learning phenomena of the NMT framework.

2 The NMT Framework

The NMT framework is proposed as a sequence-to-sequence transduction task (Sutskever et al., 2014). To make the framework clean to readers, we divide

it roughly into three independent and indispensable modules: *i*) model architecture; *ii*) model training, and *iii*) inference mechanism.

Model architecture The NMT model is usually implemented with the encoder-decoder with attention architecture. Recurrent neural networks are used to parameterize the model (Bahdanau et al., 2014; Wu et al., 2016). Then convolutional and self-attention neural networks are proposed respectively (Gehring et al., 2017; Vaswani et al., 2017).

Model training Typical training uses maximum likelihood estimation (MLE) to minimize the negative log-likelihood: $\mathcal{L}(\theta) = -\log P_{\theta}(y|x) = -\sum_t \log P_{\theta}(y_t|x, y_{<t})$. Reinforcement Learning (RL) is also leveraged for optimizing the evaluation metric based loss: $\mathcal{L}(\theta) = -\mathbb{E}_{\hat{y} \sim P_{\theta}(y|x)} m(\hat{y}, y^*)$, where $m(\cdot, y^*)$ denotes certain evaluation metric, e.g., BLEU. Like RL, minimum risk training (MRT) is also used to optimize metric (Shen et al., 2016). Besides, many tricks such as learning rate schedule, normalization techniques, and label smoothing are also used. For Non-Autoregressive neural machine Translation (NAT), knowledge distillation (KD) is used for performance boosting.

Model inference Beam search is used to find an approximate solution of the $\hat{y} = \arg_y \max P_{\theta}(y|x)$ problem. Due to several issues of the vanilla beam search, tricks like the length penalty are proposed.

3 Understanding Model Architecture

3.1 Understanding encoder/decoder

Q. Does encoder’s representation entail linguistic knowledge? Most of the works on this topic use certain linguistic tasks to assess the power of the learned hidden representations. Early on, Shi et al. (2016) begin to answer whether string-based NMT models learn about source syntax. They test hidden states’ ability to predict syntactic labels, e.g., voice, tense, smallest phrase constituent. Belinkov et al. (2017) deliver thorough analyses using the method of probing on what encoder learns about morphology knowledge of source languages. Main conclusions like the depth of the layer, the input representation (word, character), and the language types are important factors influencing the learned knowledge of morphology can be drawn from their analyses. Belinkov et al. (2020a) arrange together the analyses on the power of learned representation across various granularities of linguistic knowledge based on probing mainly for the

encoder. They further add syntactic and semantic tasks. Bisazza and Tump (2018) also study morphology knowledge captured by embedding layer.

Q. Can encoder learn word sense disambiguation? Several works attempt to understand word sense disambiguation (WSD) ability of NMT models. Rios Gonzales et al. (2017) construct a contrastive dataset where references are accompanied with a rewritten one that has an incorrect translation of a source ambiguous word. They find that the model ranks 70% of such contrastive pairs correctly, indicating the model’s strong WSD ability. Marvin and Koehn (2018) further investigate the hidden activations’ WSD ability through visualization of hidden vector clusters. Tang et al. (2019a) take the encoder as a whole for WSD ability analyses under different model architectures via probing.

Q. Does decoder’s representation entail linguistic knowledge? Belinkov et al. (2017) study the linguistic properties of decoder’s representation compared to encoder’s. They probe and find that decoder’s representation falls back a lot in accuracy of predicting POS tags. In contrast, in their later work (Belinkov et al., 2020a), they find that decoder’s representation is similar to or better than encoder’s for morphological tag prediction. Instead, Li et al. (2019a) study the possibly learned coarse-to-fine characteristics of decoder’s layer-wise representation with probes on hierarchical probing tasks.

Q. Can a single neuron entail linguistic knowledge? Instead of taking vector representation as a whole, Bau et al. (2019) leverage an unsupervised method to identify important neurons and use GMM to find neurons that controls linguistic features in prediction. Dalvi et al. (2019) also propose supervised methods to extract salient neurons and analyze their linguistic properties through probing.

Q. Can linguistic knowledge be preserved after pruning? Movva and Zhao (2020) study the representation of modules of the Transformer model while being pruned. They observe that pruning degrades semantic knowledge before affecting BLEU, and representation in higher layers changes most.

Q. Which component of NMT is more critical, encoder or decoder? Tang et al. (2019b) attempt to reveal the representational power of the encoder by removing it, so as a result, the encoder is just word and position embeddings. They find that the non-contextualized encoder representation largely

degrades performance; however, the attention module complements this as a strong feature extractor. Kasai et al. (2020) study encoder/decoder with varied depths. They find that a sufficiently deep encoder with a single-layer decoder can achieve comparable performance with balanced layer depth.

3.2 Understanding attention

3.2.1 Cross-attention

Q. Does attention learn alignments? When attention was first introduced into NMT models (Bahdanau et al., 2014), it was believed as a word alignment module inside NMT. Liu et al. (2016); Mi et al. (2016); Li et al. (2018); Baan et al. (2019) try to improve NMT’s translation quality by improving the alignment performance of its attention module. However, the attention module in NMT was far from qualified as a good word aligner compared with statistical word aligners (Koehn and Knowles, 2017; Li et al., 2019b). Although the AER of attention is dissatisfactory, Li et al. (2019b) did successfully induce decent alignments from NMT models by the method Prediction Difference (PD). Notably, Li et al. (2019b) empirically showed that, towards predictions instead of references, the performance of alignments induced by PD could surpass well-performed traditional statistical aligners. This result rekindled the confidence in inducing accurate alignments from the attention module. By improving training (Garg et al., 2019) and modeling (Alkhouli et al., 2018; Chen et al., 2020; Kobayashi et al., 2020) methods, the alignment performance of NMT’s attention are constantly improved. In the situation where translation quality is not as important as alignment performance, attention can also be extremely helpful in building a well-performed word aligner (Zenkel et al., 2020).

Q. Do attention weights reflect NMT’s reasoning? Since Bahdanau et al. (2014) introduced attention to NMT, attention weights were often claimed to explain the inner-working mechanism of neural models (Li et al., 2016). Jain and Wallace (2019) are the first to question attention’s ability to provide transparency for model predictions by showing a weak correlation between intuitive feature importance measures and attention weights in text classification, question answering, and natural language inference tasks. However, Wiegrefe and Pinter (2019) argue that Jain and Wallace (2019) does not disprove the usefulness of attention for explainability by showing the attention weights can-

not be easily hacked adversarially. Based on this observation, Moradi et al. (2020) provide a measure of the faithfulness of NMT and an adversarial regularization that can lead to more trustworthy attention heatmaps without reducing the translation quality. Current analyses are mainly focused on simpler single-head RNN based models. In the future, checking whether the current understanding holds on multi-head attention of Transformer could be an interesting direction.

3.2.2 Self-attention

Q. Is self-attention network better than RNN? The common suspicion is that self-attention can connect distant words via shorter network paths than RNNs to improve the ability to model long-range dependencies. However, this theoretical argument is not tested empirically. Tang et al. (2018) evaluate RNNs, CNNs, and multi-head attention networks (SAN) on two tasks: subject-verb agreement and word sense disambiguation to measure the ability to extract semantic features from the source text. Their experimental results show that the SAN performs distinctly better than RNNs and CNNs on word sense disambiguation. However, all of them are similar in modeling subject-verb agreement over long distances. Besides, SAN is ascribed to be weak at learning positional information of words for sequence modeling compared to the models with recurrence structure. Yang et al. (2019) show that although SAN trained on word re-ordering detection has difficulty learning positional information, SAN trained on machine translation learns better positional information than RNN.

Q. Is multi-head better than single-head? In Transformer, multi-head attention strengthens the expressive power of a model by extending a single head to multiple parallel heads. From a Bayesian perspective, An et al. (2020) understands why one needs multi-head attention by showing it is equivalent to using more samples to approximate an underlying posterior distribution. Snell et al. (2021) explain why attention obtained by MLE often correlates well with saliency and how attention can increase performance by improving its training dynamics rather than expressiveness. Raganato et al. (2020) deliver a finding that for the encoder’s multi-head attention, fixing other heads’ weight and only learning one head can achieve similar performance in high-resource translation tasks and even improve performance up to 3.5 BLEU points

in low-resource scenarios. Behnke and Heafield (2020) propose simple heuristics for pruning attention heads at the early stage of training. It confirms that most attention heads are not confident in their decisions. Michel et al. (2019); Voita et al. (2019b); Liu et al. (2021) empirically show that multi-heads are redundant at test time but are greatly helpful in training. This opens up many opportunities for downsizing these humongous models for inference.

4 Understanding Training

4.1 Training data

Q. How does data noise affect NMT? Noise in bitext corpus impacts NMT a lot. Khayrallah and Koehn (2018) investigate the impact of various types of noise of the training data on the performance of the NMT model and an SMT model. By adding many controlled types of noise to the original high-quality data, they find that the NMT model is very susceptible to noise and can degrade up to 9 BLEU points, whereas the SMT model can even obtain 1 BLUE improvement. They build five types of noise and analyze how these noises can impact translation quality. They find copy noise, where the target is just the copy of the source, is most harmful. Ott et al. (2018) reemphasize the harmfulness of copy noise in training data. They also find that beam search puts too much probability mass over the whole search space due to data uncertainty, not concentrating on accurate and relevant translations.

Q. How does the src/tgt divergence affect NMT? Briakou and Carpuat (2021) study fine-grained semantic divergences in bitext. They propose three typical divergences, lexical substitution, phrase replacement, and subtree deletion. They study their effects on NMT and find subtree deletion degrade performance the most. In a semi-supervised setting, due to extra monolingual data, the textual domains of src/tgt might exhibit topic divergence. Shen et al. (2021) propose a metric to measure such mismatch phenomenon and study its effects, particularly with varying data scales and find it can severely degrade performance in a low-resource setting.

Q. Why does DA training help? Data Augmentation (DA) methods are effective in training NMT with few theory-oriented understanding. Li et al. (2019a) borrow empirical evidence that input sensitivity and prediction margin can measure generalization ability from the learning theory community and apply them to test intrinsic changes of the

model before and after DA. DA methods generally lead to better insensitivity and a larger margin.

Q. What factors of BT data matter? Amongst all DA methods, Back-Translation (BT) is the most extensively adopted one in challenges and deployments to obtain state-of-the-art translation quality. Edunov et al. (2018a) conducts a large-scale analysis of practical BT training. They argue that *randomness* is an essential factor for improving performance, so they use sampling rather than beam search to obtain pseudo bitext. However, Caswell et al. (2019) argues that randomness might not be the reason for better practice in synthetic data generation in BT. They claim that the NMT model can automatically distinguish synthetic or real data and learn different attention patterns over them. So they propose tagged BT to improve standard BT. Following this work, Marie et al. (2020) further proves that tagged BT can prevent the NMT model from over-fitting to those machine-generated data. Besides, Graça et al. (2019) proposes a math interpretation of back-translation, which links BT to variational inference and motivates multi-turn BT.

4.2 Training loss

Q. What are the issues of NLL? *Negative log-likelihood* (NLL) loss is the default loss function to train NMT models with MLE. NLL is a token-level loss that is locally normalized and defined on ground-truth prefix. Such characteristics make NLL suffer from the following issues as discussed in Ranzato et al. (2015); Wiseman and Rush (2016): i) *exposure bias*: the model is never exposed to its own errors during training, and so the inferred histories at test-time do not resemble the gold training histories; ii) *train-test mismatch*: training uses a token-level loss, while at test-time, we target improving sequence-level evaluation metrics, such as BLEU; iii) *label bias*: the model score is locally normalized at the token level, whereas the search algorithm cares about the sequence level score. Edunov et al. (2018b) investigate other token-level loss choices such as margin-based losses and find they do not lead to significant improvement over NLL. Afterward, a large set of works have tried to propose methods based on RL to overcome the above three issues, though they seem to leave NLL unshakeable (Bengio et al., 2015; Shen et al., 2016; Wu et al., 2018; Zhang et al., 2019).

Q. Can RL-oriented loss be better than NLL? RL is used for solving pitfalls of NLL loss. How-

378 ever, large-scale experiments in Wu et al. (2016) 428
379 do not find promising performance improvements. 429
380 Later on, Wu et al. (2018) study effective training 430
381 tricks that can stably improve RL over NLL, but 431
382 analyses on why RL cannot reach our expectation is 432
383 still lacking. More recently, Choshen et al. (2020) 433
384 deliver a novel understanding of the limitations 434
385 of RL-based training. They find a peaking effect 435
386 statistics to clarify the poor exploration problem of 436
387 RL training due to the model distribution, which 437
388 renders reward for being less critical. Following 438
389 their work, Kiegeland and Kreutzer (2021) provide 439
390 several counter-evidences in terms of claims that 440
391 regard model distribution to be more critical than 441
392 reward in Choshen et al. (2020). They revisit tricks 442
393 like variance reduction, explore-exploitation trade- 443
394 off and find that peakiness cannot solely explain 444
395 improvements, and successful exploration can also 445
396 improve the likelihood of low-ranked tokens. 446

397 **Q. How does KD help with NAT?** Knowledge 428
398 Distillation at sequence-level (KD) (Kim and Rush, 429
399 2016) is another loss used to train a student NMT 430
400 model from the output distribution or prediction of 431
401 a teacher model. In Non-Autoregressive machine 432
402 Translation (NAT), KD is a crucial training tech- 433
403 nique to bring the NAT model’s performance close 434
404 to autoregressive ones (Gu et al., 2017). Zhou et al. 435
405 (2020b) investigate the critical role of KD in non- 436
406 autoregressive NMT training. They find that KD 437
407 reduces the complexity of the training bitext cor- 438
408 pora to alleviate the learning/optimization burden 439
409 of the NAT model due to its less powerful modeling 440
410 power. They also propose improved KD loss func- 441
411 tions for improved training. Xu et al. (2021) further 442
412 analyze the impacts of KD training over the intrin- 443
413 sic characteristics of the NAT model. By defining 444
414 two measures, namely word ordering agree and 445
415 lexical diversity, they empirically demonstrate that 446
416 KD is actually reducing training data complexity 447
417 in terms of word ordering and lexical choices.

418 4.3 Training tricks

419 Since Transformer has already become the de-facto 428
420 architecture for NMT best practice, several works 429
421 attempt to dig deeper into those tricks for making 430
422 Transformer training really work.

423 **Q. How does LN help?** As for the trick of Layer 428
424 Normalization (LN), Wang et al. (2019b) calculate 429
425 the instability of gradient mathematically when 430
426 putting LN layer after residual block (post-LN) and 431
427 empirically prove the effectiveness of pre-LN for

scaling up Transformer with deeper layers. Then, 428
Xiong et al. (2020) take advantage of the mean field 429
theory to prove that post-LN connection at initial- 430
ization leads to a large gradient. They find that the 431
warming-up stage is avoiding such a problem. 432

433 **Q. How residual blocks cause training instabil-** 433
434 **ity?** Besides the position of LN, Liu et al. (2020) 434
435 provide comprehensive analyses of what compli- 435
436 cates Transformer training theoretically and empir- 436
437 ically. Their analyses find that the residual blocks 437
438 can also lead to the unbalanced gradient issue. 438

439 **Q. How does label smoothing help?** As for the 439
440 typical trick label smoothing, Müller et al. (2019) 440
441 find that label smoothing can help calibrate training 441
442 instances. Gao et al. (2020) investigate its theo- 442
443 retical and empirical role. Theoretically, they find 443
444 what objective label smoothing is optimized for and 444
445 derive an analytical solution for visualization for 445
446 picking a better probability mass hyper-parameter 446
447 for smoothing (e.g., from usual 0.1 to 0.3). 447

448 5 Understanding Inference

449 5.1 Prediction explanation

450 **Q. How to attribute NMT model’s prediction?** 450
451 One effective way to interpret the NMT model’s 451
452 behavior is to understand why the model predicts 452
453 specific tokens step-wise regarding input tokens. 453
454 At the beginning of NMT, attention is leveraged to 454
455 visualize output-input correlation (Bahdanau et al., 455
456 2014). Then, Alvarez-Melis and Jaakkola (2017) 456
457 propose a perturbation-based method to collect 457
458 correlation pairs from relating every target token 458
459 to every source token so that the explanation is 459
460 model-agnostic. They exemplify with a case study 460
461 that model debugging could be conducted based 461
462 on such attention-like visualizations. Ding et al. 462
463 (2017) leverage the so-called layer-wise relevance 463
464 propagation (LRP) to capture the correlation of any 464
465 two nodes in the computation graph of the model. 465
466 They further use this method to visualize the re- 466
467 lationship between prediction and input. Several 467
468 translation errors are analyzed using LRP visual- 468
469 ization to show the power of this method. Treviso 469
470 and Martins (2020) proposes sparse/selective at- 470
471 tention as a better way than gradient and erasure 471
472 methods that relate prediction to input features (se- 472
473 quence) in terms of a success rate of a communi- 473
474 cation game (in Sec. G of the paper). Abnar and 474
475 Zuidema (2020) propose a new method for visual- 475
476 izing the flow of the information from each input 476

477	token to the output. Their proposed methods corre-	the beginning of decoding, degrading performance.	527
478	late well with attention and gradient-based method.		
479	More broadly speaking, various kinds of so-called	Q. Is beam search good enough? Stahlberg and	528
480	<i>attribution</i> methods in Sec. B.1 can all be adapted	Byrne (2019) analyze the impacts of model/search	529
481	to explain step-wise prediction of NMT. Vafa et al.	errors on performance, based on exact inference	530
482	(2021) propose a combinatorial optimization for-	for vanilla beam search. They find the model error	531
483	mulation for finding a subset of input that correlates	is more responsible. Meister et al. (2020) cleverly	532
484	most with a given output token. Their experiments	frame beam search as exact solution to a different	533
485	show that the proposed method is most faithful	decoding objective to gain insights into why high	534
486	among other explanation methods.	probability under a model alone may not guaran-	535
487		tee adequacy. Eikema and Aziz (2020) attempt to	536
488	Q. How to properly evaluate prediction attribu-	clarify the problem of maximum a posterior (MAP)	537
489	tion of NMT model? As mentioned in the previ-	based beam search. They find that translation distri-	538
490	ous question, various prediction attribution meth-	butions of the model do reproduce various statistics	539
491	ods can be used to explain model prediction. How-	of the training data, but beam search strays from	540
492	ever, in practice which method to choose? There	such statistics. They also propose to use Minimum	541
493	seems to be no fixed answer currently since there	Bayes Risk (MBR) decoding instead. Müller and	542
494	are already several issues found with attributions.	Sennrich (2021) study the properties of MBR de-	543
495	Regardless of these issues, several works have pro-	coding. They find that MBR decoding still exhibits	544
496	posed methods to evaluate attributions from differ-	length and token frequency bias due to the bias of	545
497	ent perspectives. Li et al. (2020) propose to use the	evaluation metrics, but MBR also increases robust-	546
498	word-to-word correlation rules extracted by vari-	ness against copy noise and domain shift.	547
499	ous attribution methods to train models close to the		
500	original NMT model and uses the closeness as a	6 Understanding Model Behavior	548
501	way to evaluate the attribution results. Treviso and		
502	Martins (2020)’s communication game can also	Model behavior understanding is generally an ef-	549
503	be used as an evaluation method. Beyond NMT,	fort to <i>characterize</i> and <i>analyze</i> certain property of	550
504	several works propose methods and benchmarks	the <i>model’s predictions</i> in terms of certain aspect or	551
505	for evaluating attribution (Hao, 2020; Arras et al.,	factor in concern. Currently, we can divide those re-	552
506	2019; Ismail et al., 2020; Ding and Koehn, 2021).	search works into <u>three main categories</u> according	553
507		to their adopted analysis methodologies: a). static	554
508	5.2 Decoding explanation	analysis: that tries to <i>directly</i> analyze properties	555
509		of the model’s predictions, e.g., fluency, grammat-	556
510	Q. Do larger beams lead to better results? Af-	icality, word choice, the degree of literalness or	557
511	ter widely adoption of NMT, Koehn and Knowles	creativity, etc.; b). controlled analysis: that tries	558
512	(2017) describe a common phenomenon of beam	to characterize the model’s reaction to inputs <i>con-</i>	559
513	search decoding across various language pairs, that	<i>structed</i> with certain properties in concern, e.g.,	560
514	is, by increasing beam size, the BLEU score will	compositionality, specific linguistic phenomenon,	561
515	rise up a little and then jump down quite a lot,	etc.; c). dynamic analysis: that attempts to do <i>in-</i>	562
516	compared to SMT. This represents the so-called	<i>terventions</i> and <i>manipulation</i> to the inputs or the	563
517	<i>length bias</i> problem which has been investigated to	model, which might help reveal weaknesses of the	564
518	show its correlation with i) decoding scoring func-	model when making predictions about these inputs,	565
519	tion (Huang et al., 2018; Yang et al., 2018) and ii)	e.g., adversarial examples, syntactic/semantic vari-	566
520	beam size. Murray and Chiang (2018) further find	ants, hallucinations, noise in training data, etc.	567
521	that label bias is one factor of such a problem and		
522	propose a simple heuristic to alleviate it. Cohen	6.1 Static analysis	568
523	and Beck (2019) deliver a more detailed analysis		
524	on beam search using the concept of search discrep-	Q. Is NMT model’s prediction linguistically nat-	569
525	ancies, which is computed through the difference	ural? Toral and Sánchez-Cartagena (2017) con-	570
526	between the maximum log probability and log prob-	duct a comparative study between the predictions	571
	ability of the ground-truth token at every time-step	of NMT and SMT models in terms of fluency, re-	572
	under force decoding. They find that a larger beam	ordering, sentence length among 9 language direc-	573
	size may cause larger and more discrepancies at	tions. It highlights the power of NMT models on	574

generating more fluent, accurately reordered predictions. Later on, [Martindale et al. \(2019\)](#) also study the fluency-adequacy dilemma of neural models. [Wei et al. \(2018\)](#) investigate the grammaticality of NMT outputs. They leverage the so-called English Resource Grammar as a reference for comparison. They find that over 93% of the model translations are parseable, suggesting that the model learns to generate conforming to grammar; however, rare syntactic rules are seldom learned.

Q. Can NMT model generate long-tailed translation? Long-tailed translations can be predictions that contains low-frequent tokens, complex phrases, and advanced sentence structures. [Rau-nak et al. \(2020\)](#) characterize the hardness of NMT to predict long-tailed words and tokens through token-level and sentence-level metrics. [Agrawal and Carpuat \(2019\)](#) study text complexity of predictions and focus on controlling the outputs towards less complexity. [Vanmassenhove et al. \(2021\)](#) give a detailed and sufficient analysis on the richness of word choices and synonyms etc.. They also design several metrics to evaluate linguistic complexity. Long-tailed translations can also be indirect translations of phrases that are seldom in the common bitext corpus. [Zhai et al. \(2020\)](#) investigate whether NMT models are capable of producing non-literal translations. They propose methods to detect those non-literal translation phenomena in bitext.

Q. Can model’s prediction be well calibrated? Calibration is a sound property of a learned model to predict the probability of the true correctness likelihood ([Guo et al., 2017](#)). [Kumar and Sarawagi \(2019\)](#) analyze the sources of surprising miscalibration in NMT. They find that the severe miscalibration of the EOS token and the suppression of attention uncertainty are two main reasons. [Wang et al. \(2020\)](#) further study the fine-grained calibration of the model predictions. They characterize miscalibrated tokens with linguistic features, such as questions about how part-of-speech, frequency, word position, word granularity affect calibration.

6.2 Controlled analysis

Q. Can NMT model handle inputs with different types of linguistic phenomenon? Inputs to an NMT model can be linguistically sophisticated. [Burchardt et al. \(2017\)](#) manually construct a test suite with different kinds of linguistic phenomenon of the source input sentences, for instances, multi-word expressions, verb tense/aspect/mood, named

entity, and terminology in German \leftrightarrow English translation tasks. They compare the performance of the Google NMT system at that time with SMT and rule-based models on this test suite. They find that neural models handle multiword expressions much better than rule-based and SMT models, while rule-based ones can handle verb tense/aspect/mood structure the best, and SMT handles named entities the best. Similarly, [Isabelle et al. \(2017\)](#) construct a challenge set with yes/no questions for analyzing both phrase and neural translation models’ capability to handle three categories of linguistic phenomenon in English \Rightarrow French task. They find NMT models are much better at tackling subject-verb agreement and perform well on handling both lexico-syntactic and syntactic divergences. They also identify some weaknesses of neural models; please refer to Table 3 in that paper for details.

Q. Can NMT model handle inputs compositionally? [Raunak et al. \(2019\)](#) measure two distinct traits of compositionality - *productivity* and *systematicity* - of the NMT model by comparing performance before and after sentence concatenation. Their experiments quantitatively attribute the poor performance to the weakness of the encoder’s representational power. [Li et al. \(2021\)](#) build a benchmark for training and testing the model’s compositional capability to tackle *compounds*, which are constructed through pre-defined atoms and syntactic rules. [Dankers et al. \(2021\)](#) evaluate the model’s compositionality through the lens of the model’s local/global processing of the input. [Voita et al. \(2019a\)](#) focus on a problem of NMT model trained on sentence-level, that is, while the model can accurately translate sentences A and B, but can not when A and B are concatenated in a broader context, which can be also regarded as a problem of compositionality in discourse translation. All the above works find that Transformer or more or other NMT models have poor compositionality.

6.3 Dynamic analysis

Q. Is NMT model robust to inputs? Adversarial examples are an essential direction for testing the NMT model’s robustness where the adversarial inputs are created through input manipulation. [Belinkov and Bisk \(2018\)](#) is the first to investigate how realistic, natural adversarial input (e.g. character-level keyboard typing errors) can break the char-based translation model. [Zhao et al. \(2017\)](#) and [Cheng et al. \(2020\)](#) investigate model-based

675 methods for generating adversarial examples; while
 676 [Ebrahimi et al. \(2018\)](#) focus on attacking the char-
 677 based NMT models. Besides adversarials, [He et al.](#)
 678 [\(2019\)](#) use different input word replacement strate-
 679 gies to identify important source words that guaran-
 680 tee the translation quality of the source input. They
 681 argue that the so-called *importance words* are cru-
 682 cial to guarantee fertility and should not be ignored.
 683 [Fadaee and Monz \(2020\)](#) study the so-called *volatil-*
 684 *ity* of NMT models where the input is semantically
 685 and syntactically transformed while the prediction
 686 can have unexpected disastrous changes. They
 687 find RNN and Transformer display volatile behav-
 688 ior in 26% and 19% of sentence variations.

689 **Q. When or why does NMT model hallucinate?**

690 Hallucination is a recently identified phenomenon
 691 in [Lee et al. \(2018\)](#). It is the problem of an NMT
 692 model that outputs irrelevant sentence predictions
 693 or textual spans with respect to certain *constructed*
 694 input. They analyze the attention patterns that dis-
 695 tinguish hallucinated and normal predictions. [Rau-](#)
 696 [nak et al. \(2021\)](#) connect this phenomenon to long-
 697 tailed memorization effect of the model. [Wang and](#)
 698 [Sennrich \(2020\)](#) regard exposure bias as one factor
 699 of hallucination and find domain-shift amplifies its
 700 harmfulness. [Zhou et al. \(2020a\)](#) tackle the identifi-
 701 cation problem of hallucination of neural sequence
 702 model in general. They construct datasets for token-
 703 wise annotation of hallucination and explore some
 704 basic methods for detecting hallucinated tokens.

705 **7 Limitations, Future and Conclusion**

706 In this part, we summarize several current limita-
 707 tions of those aforementioned understandings, in-
 708 terpretations and findings, and propose a few future
 709 directions on the understanding course of NMT.

- 710 • Vacuousness of representation probing: prob-
 711 ing measures the feature generalization ability
 712 of the NMT learned representations on certain
 713 concerned linguistic task, however, does do-
 714 ing well on that task really help the model with
 715 the translation task? Such direct correlation
 716 between probing task and translation is very
 717 vague as well. [Elazar et al. \(2021\)](#) attempt
 718 to resolve this issue through explicit remov-
 719 ing certain linguistic knowledge in the learned
 720 representation of BERT to see its utility on
 721 the downstream classification tasks. So, how
 722 about using such analysis in more complex
 723 translation tasks ([Ravichander et al., 2021](#)).

- 724 • Usability of prediction attribution: § 5’s first
 725 question discusses many methods for attribut-
 726 ing predicted tokens to previous input tokens.
 727 Besides the evaluation issue of these methods,
 728 how to use such attributions to debug model,
 729 moreover, to improve user trust beyond sole
 730 alignment or to improve interactive transla-
 731 tion ([Santy et al., 2019](#)) is not well explored.
- 732 • Insufficient understanding on learning dynam-
 733 ics: by exploring learning dynamics, theorists
 734 have found critical learning phases that deter-
 735 mine final generalization ([Achille et al., 2017](#);
 736 [Hu et al., 2020](#); [Jastrzebski et al., 2021](#)). How-
 737 ever, investigations of learning dynamics are
 738 largely neglected in NMT, except for [Saphra](#)
 739 [and Lopez \(2019\)](#); [Zhu et al. \(2020\)](#); [Voita](#)
 740 [et al. \(2021\)](#). We think gaining more insights
 741 in the learning dynamics of NMT model might
 742 help with better curriculum, data selection, in-
 743 stance reweighting, noise-based learning, etc..
- 744 • Lack of data-centric understanding: many of
 745 the current understandings leverage a model-
 746 centric analysis, i.e., only considering archi-
 747 tectural inductive bias without knowing char-
 748 acteristics of the training data, however, the
 749 ultimate model behavior is largely determined
 750 by the training instances as well ([Yona et al.,](#)
 751 [2021](#)). In NLP, there have been works that
 752 using dataset attribution techniques like in-
 753 fluence function ([Koh and Liang, 2017](#)) to
 754 find artifacts in the training set for text clas-
 755 sification ([Han et al., 2020](#)). Thus how to
 756 adopt similar methods to the complex ma-
 757 chine translation task should be studied. We
 758 think this direction may help researchers cu-
 759 rate more compact and continuously-updated
 760 datasets for sample-efficient training and con-
 761 tinual learning ([Cao et al., 2021](#)) of NMT.

762 As a conclusion, the understanding of the evol-
 763 ving NMT framework should be always on its way
 764 and, to find limitations of the current best prac-
 765 tice, emerging topics with multilingual, continual
 766 and discourse NMT ([Dabre et al., 2020](#); [Garcia](#)
 767 [et al., 2021](#); [Yin et al., 2021](#)) require better un-
 768 derstanding, theory-oriented and empirical analyses as
 769 well, so the FAQs here ([https://nmtology.](https://nmtology.github.io/)
 770 [github.io/](https://nmtology.github.io/)) might and should be revisited and
 771 updated in new scenarios. The authors believe that
 772 knowing the historic understandings could help the
 773 community pave the way towards the future.

774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830

References

Samira Abnar and Willem Zuidema. 2020. [Quantifying attention flow in transformers](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4190–4197, Online. Association for Computational Linguistics.

Alessandro Achille, Matteo Rovere, and Stefano Soatto. 2017. Critical learning periods in deep neural networks. *arXiv preprint arXiv:1711.08856*.

Sweta Agrawal and Marine Carpuat. 2019. [Controlling text complexity in neural machine translation](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 1549–1564, Hong Kong, China. Association for Computational Linguistics.

Guillaume Alain and Yoshua Bengio. 2016. Understanding intermediate layers using linear classifier probes. *arXiv preprint arXiv:1610.01644*.

Tamer Alkhouli, Gabriel Bretschner, and Hermann Ney. 2018. [On the alignment problem in multi-head attention-based neural machine translation](#). In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 177–185, Brussels, Belgium. Association for Computational Linguistics.

David Alvarez-Melis and Tommi Jaakkola. 2017. [A causal framework for explaining the predictions of black-box sequence-to-sequence models](#). In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 412–421, Copenhagen, Denmark. Association for Computational Linguistics.

Bang An, Jie Lyu, Zhenyi Wang, Chunyuan Li, Changwei Hu, Fei Tan, Ruiyi Zhang, Yifan Hu, and Changyou Chen. 2020. [Repulsive attention: Rethinking multi-head attention as Bayesian inference](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 236–255, Online. Association for Computational Linguistics.

Leila Arras, Ahmed Osman, Klaus-Robert Müller, and Wojciech Samek. 2019. [Evaluating recurrent neural network explanations](#). In *Proceedings of the 2019 ACL Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 113–126, Florence, Italy. Association for Computational Linguistics.

Joris Baan, Jana Leible, Mitja Nikolaus, David Rau, Dennis Ulmer, Tim Baumgärtner, Dieuwke Hupkes, and Elia Bruni. 2019. [On the realization of compositionality in neural networks](#). In *Proceedings of the 2019 ACL Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 127–137, Florence, Italy. Association for Computational Linguistics.

Sebastian Bach, Alexander Binder, Grégoire Montavon, Frederick Klauschen, Klaus-Robert Müller, and Wojciech Samek. 2015. On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation. *PLoS one*, 10(7):e0130140. 831
832
833
834
835

Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2014. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*. 836
837
838
839

Loïc Barrault, Magdalena Biesialska, Ondřej Bojar, Marta R Costa-jussà, Christian Federmann, Yvette Graham, Roman Grundkiewicz, Barry Haddow, Matthias Huck, Eric Joanis, et al. 2020. Findings of the 2020 conference on machine translation (wmt20). In *Proceedings of the Fifth Conference on Machine Translation*, pages 1–55. 840
841
842
843
844
845
846

Loïc Barrault, Ondřej Bojar, Marta R Costa-Jussa, Christian Federmann, Mark Fishel, Yvette Graham, Barry Haddow, Matthias Huck, Philipp Koehn, Shervin Malmasi, et al. 2019. Findings of the 2019 conference on machine translation (wmt19). In *Proceedings of the Fourth Conference on Machine Translation (Volume 2: Shared Task Papers, Day 1)*, pages 1–61. 847
848
849
850
851
852
853
854

Osbert Bastani, Carolyn Kim, and Hamsa Bastani. 2017. Interpreting blackbox models via model extraction. *arXiv preprint arXiv:1705.08504*. 855
856
857

Anthony Bau, Yonatan Belinkov, Hassan Sajjad, Nadir Durrani, Fahim Dalvi, and James R. Glass. 2019. [Identifying and controlling important neurons in neural machine translation](#). In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net. 858
859
860
861
862
863
864

Maximiliana Behnke and Kenneth Heafield. 2020. [Losing heads in the lottery: Pruning transformer attention in neural machine translation](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2664–2674, Online. Association for Computational Linguistics. 865
866
867
868
869
870
871

Yonatan Belinkov and Yonatan Bisk. 2018. [Synthetic and natural noise both break neural machine translation](#). In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net. 872
873
874
875
876
877

Yonatan Belinkov, Nadir Durrani, Fahim Dalvi, Hassan Sajjad, and James Glass. 2017. [What do neural machine translation models learn about morphology?](#) In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 861–872, Vancouver, Canada. Association for Computational Linguistics. 878
879
880
881
882
883
884

Yonatan Belinkov, Nadir Durrani, Fahim Dalvi, Hassan Sajjad, and James Glass. 2020a. [On the linguistic](#) 885
886

887	representational power of neural machine translation	Yun Chen, Yang Liu, Guanhua Chen, Xin Jiang, and	941
888	models. <i>Computational Linguistics</i> , 46(1):1–52.	Qun Liu. 2020. Accurate word alignment induction	942
889	Yonatan Belinkov, Sebastian Gehrmann, and Ellie	from neural machine translation. In <i>Proceedings of</i>	943
890	Pavlick. 2020b. Interpretability and analysis in neural	<i>the 2020 Conference on Empirical Methods in Natural</i>	944
891	NLP. In <i>Proceedings of the 58th Annual Meeting of the</i>	<i>Language Processing (EMNLP)</i> , pages 566–576,	945
892	<i>Association for Computational Linguistics: Tutorial Abstracts</i> ,	Online. Association for Computational Linguistics.	946
893	pages 1–5, Online. Association for Computational Linguistics.		
894			
895	Yonatan Belinkov and James Glass. 2019. Analysis	Minhao Cheng, Jinfeng Yi, Pin-Yu Chen, Huan Zhang,	947
896	methods in neural language processing: A survey.	and Cho-Jui Hsieh. 2020. Seq2sick: Evaluating the	948
897	<i>Transactions of the Association for Computational</i>	robustness of sequence-to-sequence models with ad-	949
898	<i>Linguistics</i> , 7:49–72.	versarial examples. In <i>The Thirty-Fourth AAAI Con-</i>	950
899		<i>ference on Artificial Intelligence, AAAI 2020, The</i>	951
900		<i>Thirty-Second Innovative Applications of Artificial</i>	952
901		<i>Intelligence Conference, IAAI 2020, The Tenth AAAI</i>	953
902		<i>Symposium on Educational Advances in Artificial In-</i>	954
903		<i>telligence, EAAI 2020, New York, NY, USA, Febru-</i>	955
904		<i>ary 7-12, 2020</i> , pages 3601–3608. AAAI Press.	956
905			
906		David Chiang. 2005. A hierarchical phrase-based	957
907		model for statistical machine translation. In <i>Pro-</i>	958
908		<i>ceedings of the 43rd annual meeting of the associ-</i>	959
909		<i>ation for computational linguistics (acl’05)</i> , pages	960
910		263–270.	961
911			
912		Leshem Choshen, Lior Fox, Zohar Aizenbud, and Omri	962
913		Abend. 2020. On the weaknesses of reinforcement	963
914		learning for neural machine translation. In <i>8th Inter-</i>	964
915		<i>national Conference on Learning Representations,</i>	965
916		<i>ICLR 2020, Addis Ababa, Ethiopia, April 26-30,</i>	966
917		<i>2020</i> . OpenReview.net.	967
918			
919		Eldan Cohen and Christopher Beck. 2019. Empirical	968
920		analysis of beam search performance degradation in	969
921		neural sequence models. In <i>Proceedings of the 36th</i>	970
922		<i>International Conference on Machine Learning</i> , vol-	971
923		ume 97 of <i>Proceedings of Machine Learning Re-</i>	972
924		<i>search</i> , pages 1290–1299. PMLR.	973
925			
926		Alexis Conneau, German Kruszewski, Guillaume Lam-	974
927		ple, Loïc Barrault, and Marco Baroni. 2018. What	975
928		you can cram into a single \$&!#* vector: Probing	976
929		sentence embeddings for linguistic properties. In	977
930		<i>Proceedings of the 56th Annual Meeting of the As-</i>	978
931		<i>sociation for Computational Linguistics (Volume 1:</i>	979
932		<i>Long Papers)</i> , pages 2126–2136, Melbourne, Aus-	980
933		tralia. Association for Computational Linguistics.	981
934			
935			
936		Raj Dabre, Chenhui Chu, and Anoop Kunchukut-	982
937		tan. 2020. A survey of multilingual neural ma-	983
938		chine translation. <i>ACM Computing Surveys (CSUR)</i> ,	984
939		53(5):1–38.	985
940			
941		Fahim Dalvi, Nadir Durrani, Hassan Sajjad, Yonatan	986
942		Belinkov, Anthony Bau, and James R. Glass. 2019.	987
943		What is one grain of sand in the desert? analyzing	988
944		individual neurons in deep NLP models. In <i>The</i>	989
945		<i>Thirty-Third AAAI Conference on Artificial Intelli-</i>	990
946		<i>gence, AAAI 2019, The Thirty-First Innovative Ap-</i>	991
947		<i>plications of Artificial Intelligence Conference, IAAI</i>	992
948		<i>2019, The Ninth AAAI Symposium on Educational</i>	993
949		<i>Advances in Artificial Intelligence, EAAI 2019, Hon-</i>	994
950		<i>olulu, Hawaii, USA, January 27 - February 1, 2019,</i>	995
951		pages 6309–6317. AAAI Press.	996
952			
953		Marina Danilevsky, Kun Qian, Ranit Aharonov, Yan-	997
954		nis Katsis, Ban Kawas, and Prithviraj Sen. 2020. A	998

999	survey of the state of explainable AI for natural language processing. In <i>Proceedings of the 1st Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 10th International Joint Conference on Natural Language Processing</i> , pages 447–459, Suzhou, China. Association for Computational Linguistics.	1054
1000		1055
1001		1056
1002		1057
1003		1058
1004		
1005		
1006	Verna Dankers, Elia Bruni, and Dieuwke Hupkes. 2021. The paradox of the compositionality of natural language: a neural machine translation case study. <i>arXiv preprint arXiv:2108.05885</i> .	1059
1007		1060
1008		1061
1009		1062
1010	Shuoyang Ding and Philipp Koehn. 2021. Evaluating saliency methods for neural language models. <i>arXiv preprint arXiv:2104.05824</i> .	1063
1011		1064
1012		1065
1013	Shuoyang Ding, Hainan Xu, and Philipp Koehn. 2019. Saliency-driven word alignment interpretation for neural machine translation. In <i>Proceedings of the Fourth Conference on Machine Translation (Volume 1: Research Papers)</i> , pages 1–12, Florence, Italy. Association for Computational Linguistics.	1066
1014		1067
1015		1068
1016		1069
1017		1070
1018		1071
1019	Yanzhuo Ding, Yang Liu, Huanbo Luan, and Maosong Sun. 2017. Visualizing and understanding neural machine translation. In <i>Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 1150–1159, Vancouver, Canada. Association for Computational Linguistics.	1072
1020		1073
1021		1074
1022		1075
1023		1076
1024		1077
1025		1078
1026	Javid Ebrahimi, Daniel Lowd, and Dejing Dou. 2018. On adversarial examples for character-level neural machine translation. In <i>Proceedings of the 27th International Conference on Computational Linguistics</i> , pages 653–663, Santa Fe, New Mexico, USA. Association for Computational Linguistics.	1079
1027		1080
1028		1081
1029		1082
1030		1083
1031		1084
1032	Sergey Edunov, Myle Ott, Michael Auli, and David Grangier. 2018a. Understanding back-translation at scale. In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 489–500, Brussels, Belgium. Association for Computational Linguistics.	1085
1033		1086
1034		1087
1035		1088
1036		1089
1037		1090
1038	Sergey Edunov, Myle Ott, Michael Auli, David Grangier, and Marc’Aurelio Ranzato. 2018b. Classical structured prediction losses for sequence to sequence learning. In <i>Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)</i> , pages 355–364, New Orleans, Louisiana. Association for Computational Linguistics.	1091
1039		1092
1040		1093
1041		1094
1042		1095
1043		1096
1044		1097
1045		1098
1046		1099
1047	Bryan Eikema and Wilker Aziz. 2020. Is MAP decoding all you need? the inadequacy of the mode in neural machine translation. In <i>Proceedings of the 28th International Conference on Computational Linguistics</i> , pages 4506–4520, Barcelona, Spain (Online). International Committee on Computational Linguistics.	1100
1048		1101
1049		1102
1050		1103
1051		1104
1052		1105
1053		1106
	Yanai Elazar, Shauli Ravfogel, Alon Jacovi, and Yoav Goldberg. 2021. Amnesic Probing: Behavioral Explanation with Amnesic Counterfactuals. <i>Transactions of the Association for Computational Linguistics</i> , 9:160–175.	1107
		1108
	Marzieh Fadaee and Christof Monz. 2020. The unreasonable volatility of neural machine translation models. In <i>Proceedings of the Fourth Workshop on Neural Generation and Translation</i> , pages 88–96, Online. Association for Computational Linguistics.	1109
		1110
	Yingbo Gao, Weiyue Wang, Christian Herold, Zijian Yang, and Hermann Ney. 2020. Towards a better understanding of label smoothing in neural machine translation. In <i>Proceedings of the 1st Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 10th International Joint Conference on Natural Language Processing</i> , pages 212–223, Suzhou, China. Association for Computational Linguistics.	1111
		1112
	Xavier Garcia, Noah Constant, Ankur Parikh, and Orhan Firat. 2021. Towards continual learning for multilingual machine translation via vocabulary substitution. In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 1184–1192, Online. Association for Computational Linguistics.	1113
		1114
	Sarthak Garg, Stephan Peitz, Udhyakumar Nallasamy, and Matthias Paulik. 2019. Jointly learning to align and translate with transformer models. In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 4453–4462, Hong Kong, China. Association for Computational Linguistics.	1115
		1116
	Jon Gauthier, Jennifer Hu, Ethan Wilcox, Peng Qian, and Roger Levy. 2020. SyntaxGym: An online platform for targeted evaluation of language models. In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations</i> , pages 70–76, Online. Association for Computational Linguistics.	1117
		1118
	Jonas Gehring, Michael Auli, David Grangier, Denis Yarats, and Yann N Dauphin. 2017. Convolutional sequence to sequence learning. In <i>International Conference on Machine Learning</i> , pages 1243–1252. PMLR.	1119
		1120
	Miguel Graça, Yunsu Kim, Julian Schamper, Shahram Khadivi, and Hermann Ney. 2019. Generalizing back-translation in neural machine translation. In <i>Proceedings of the Fourth Conference on Machine Translation (Volume 1: Research Papers)</i> , pages 45–52, Florence, Italy. Association for Computational Linguistics.	1121
		1122
	Filip Graliński, Anna Wróblewska, Tomasz Stańsławek, Kamil Grabowski, and Tomasz	1123

1111	Górecki. 2019. GEval: Tool for debugging NLP datasets and models . In <i>Proceedings of the 2019 ACL Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP</i> , pages 254–262, Florence, Italy. Association for Computational Linguistics.	1166
1112		1167
1113		1168
1114		1169
1115		1170
1116		1171
1117	Jiatao Gu, James Bradbury, Caiming Xiong, Victor OK Li, and Richard Socher. 2017. Non-autoregressive neural machine translation. <i>arXiv preprint arXiv:1711.02281</i> .	1172
1118		1173
1119		1174
1120		1175
1121	Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q Weinberger. 2017. On calibration of modern neural networks. In <i>International Conference on Machine Learning</i> , pages 1321–1330. PMLR.	1176
1122		1177
1123		1178
1124		1179
1125	Xiaochuang Han, Byron C. Wallace, and Yulia Tsvetkov. 2020. Explaining black box predictions and unveiling data artifacts through influence functions . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 5553–5563, Online. Association for Computational Linguistics.	1180
1126		1181
1127		1182
1128		1183
1129		1184
1130		1185
1131		1186
1132	Yiding Hao. 2020. Evaluating attribution methods using white-box LSTMs . In <i>Proceedings of the Third BlackboxNLP Workshop on Analyzing and Interpreting Neural Networks for NLP</i> , pages 300–313, Online. Association for Computational Linguistics.	1187
1133		1188
1134		1189
1135		1190
1136		1191
1137	Shilin He, Zhaopeng Tu, Xing Wang, Longyue Wang, Michael Lyu, and Shuming Shi. 2019. Towards understanding neural machine translation with word importance . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 953–962, Hong Kong, China. Association for Computational Linguistics.	1192
1138		1193
1139		1194
1140		1195
1141		1196
1142		1197
1143		1198
1144		1199
1145		1200
1146	Wei Hu, Lechao Xiao, Ben Adlam, and Jeffrey Pennington. 2020. The surprising simplicity of the early-time learning dynamics of neural networks. <i>arXiv preprint arXiv:2006.14599</i> .	1201
1147		1202
1148		1203
1149		1204
1150	Liang Huang, Kai Zhao, and Mingbo Ma. 2018. When to finish? optimal beam search for neural text generation (modulo beam size). <i>arXiv preprint arXiv:1809.00069</i> .	1205
1151		1206
1152		1207
1153		1208
1154	Pierre Isabelle, Colin Cherry, and George Foster. 2017. A challenge set approach to evaluating machine translation . In <i>Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing</i> , pages 2486–2496, Copenhagen, Denmark. Association for Computational Linguistics.	1209
1155		1210
1156		1211
1157		1212
1158		1213
1159		1214
1160	Aya Abdelsalam Ismail, Mohamed Gunady, Hector Corrada Bravo, and Soheil Feizi. 2020. Benchmarking deep learning interpretability in time series predictions . In <i>Advances in Neural Information Processing Systems</i> , volume 33, pages 6441–6452. Curran Associates, Inc.	1215
1161		1216
1162		1217
1163		1218
1164		1219
1165		1220
	Sarthak Jain and Byron C. Wallace. 2019. Attention is not Explanation . In <i>Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)</i> , pages 3543–3556, Minneapolis, Minnesota. Association for Computational Linguistics.	1221
		1222
	Stanislaw Jastrzebski, Devansh Arpit, Oliver Astrand, Giancarlo B Kerg, Huan Wang, Caiming Xiong, Richard Socher, Kyunghyun Cho, and Krzysztof J Geras. 2021. Catastrophic fisher explosion: Early phase fisher matrix impacts generalization. In <i>International Conference on Machine Learning</i> , pages 4772–4784. PMLR.	
	Jungo Kasai, Nikolaos Pappas, Hao Peng, James Cross, and Noah A Smith. 2020. Deep encoder, shallow decoder: Reevaluating non-autoregressive machine translation. <i>arXiv preprint arXiv:2006.10369</i> .	
	Huda Khayrallah and Philipp Koehn. 2018. On the impact of various types of noise on neural machine translation . In <i>Proceedings of the 2nd Workshop on Neural Machine Translation and Generation</i> , pages 74–83, Melbourne, Australia. Association for Computational Linguistics.	
	Samuel Kiegeand and Julia Kreutzer. 2021. Revisiting the weaknesses of reinforcement learning for neural machine translation . In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 1673–1681, Online. Association for Computational Linguistics.	
	Yoon Kim and Alexander M. Rush. 2016. Sequence-level knowledge distillation . In <i>Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing</i> , pages 1317–1327, Austin, Texas. Association for Computational Linguistics.	
	Goro Kobayashi, Tatsuki Kuribayashi, Sho Yokoi, and Kentaro Inui. 2020. Attention is not only a weight: Analyzing transformers with vector norms . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 7057–7075, Online. Association for Computational Linguistics.	
	Philipp Koehn and Rebecca Knowles. 2017. Six challenges for neural machine translation . In <i>Proceedings of the First Workshop on Neural Machine Translation</i> , pages 28–39, Vancouver. Association for Computational Linguistics.	
	Philipp Koehn, Franz J. Och, and Daniel Marcu. 2003. Statistical phrase-based translation . In <i>Proceedings of the 2003 Human Language Technology Conference of the North American Chapter of the Association for Computational Linguistics</i> , pages 127–133.	
	Pang Wei Koh and Percy Liang. 2017. Understanding black-box predictions via influence functions. In <i>International Conference on Machine Learning</i> , pages 1885–1894. PMLR.	

1223	Aviral Kumar and Sunita Sarawagi. 2019. Calibration of encoder decoder models for neural machine translation. <i>arXiv preprint arXiv:1903.00802</i> .	
1224		
1225		
1226	Katherine Lee, Orhan Firat, Ashish Agarwal, Clara Fannjiang, and David Sussillo. 2018. Hallucinations in neural machine translation. <i>OpenReview.net</i> .	
1227		
1228		
1229	Piyawat Lertvittayakumjorn and Francesca Toni. 2021. Explanation-based human debugging of nlp models: A survey. <i>arXiv preprint arXiv:2104.15135</i> .	
1230		
1231		
1232	Guanlin Li, Lemao Liu, Xintong Li, Conghui Zhu, Tiejun Zhao, and Shuming Shi. 2019a. Understanding and Improving Hidden Representations for Neural Machine Translation . In <i>Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)</i> , pages 466–477, Minneapolis, Minnesota. Association for Computational Linguistics.	
1233		
1234		
1235		
1236		
1237		
1238		
1239		
1240		
1241	Jierui Li, Lemao Liu, Huayang Li, Guanlin Li, Guoping Huang, and Shuming Shi. 2020. Evaluating explanation methods for neural machine translation . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 365–375, Online. Association for Computational Linguistics.	
1242		
1243		
1244		
1245		
1246		
1247		
1248	Jiwei Li, Will Monroe, and Dan Jurafsky. 2016. Understanding neural networks through representation erasure. <i>arXiv preprint arXiv:1612.08220</i> .	
1249		
1250		
1251	Xintong Li, Guanlin Li, Lemao Liu, Max Meng, and Shuming Shi. 2019b. On the word alignment from neural machine translation . In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 1293–1303, Florence, Italy. Association for Computational Linguistics.	
1252		
1253		
1254		
1255		
1256		
1257	Xintong Li, Lemao Liu, Zhaopeng Tu, Shuming Shi, and Max Meng. 2018. Target foresight based attention for neural machine translation . In <i>Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)</i> , pages 1380–1390, New Orleans, Louisiana. Association for Computational Linguistics.	
1258		
1259		
1260		
1261		
1262		
1263		
1264		
1265	Yafu Li, Yongjing Yin, Yulong Chen, and Yue Zhang. 2021. On compositional generalization of neural machine translation . In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 4767–4780, Online. Association for Computational Linguistics.	
1266		
1267		
1268		
1269		
1270		
1271		
1272		
1273	Lemao Liu, Masao Utiyama, Andrew Finch, and Ei-ichiro Sumita. 2016. Neural machine translation with supervised attention . In <i>Proceedings of COLING 2016, the 26th International Conference on Computational Linguistics: Technical Papers</i> , pages 3093–3102, Osaka, Japan. The COLING 2016 Organizing Committee.	
1274		
1275		
1276		
1277		
1278		
1279		
	Liyuan Liu, Jialu Liu, and Jiawei Han. 2021. Multi-head or single-head? an empirical comparison for transformer training. <i>arXiv preprint arXiv:2106.09650</i> .	1280
		1281
		1282
		1283
	Liyuan Liu, Xiaodong Liu, Jianfeng Gao, Weizhu Chen, and Jiawei Han. 2020. Understanding the difficulty of training transformers . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 5747–5763, Online. Association for Computational Linguistics.	1284
		1285
		1286
		1287
		1288
		1289
		1290
	Siwen Luo, Hamish Ivison, Caren Han, and Josiah Poon. 2021. Local interpretations for explainable natural language processing: A survey. <i>arXiv preprint arXiv:2103.11072</i> .	1291
		1292
		1293
		1294
	Benjamin Marie, Raphael Rubino, and Atsushi Fujita. 2020. Tagged back-translation revisited: Why does it really work? In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 5990–5997, Online. Association for Computational Linguistics.	1295
		1296
		1297
		1298
		1299
		1300
	Marianna Martindale, Marine Carpuat, Kevin Duh, and Paul McNamee. 2019. Identifying fluently inadequate output in neural and statistical machine translation . In <i>Proceedings of Machine Translation Summit XVII: Research Track</i> , pages 233–243, Dublin, Ireland. European Association for Machine Translation.	1301
		1302
		1303
		1304
		1305
		1306
		1307
	Rebecca Marvin and Philipp Koehn. 2018. Exploring word sense disambiguation abilities of neural machine translation systems (non-archival extended abstract) . In <i>Proceedings of the 13th Conference of the Association for Machine Translation in the Americas (Volume 1: Research Track)</i> , pages 125–131, Boston, MA. Association for Machine Translation in the Americas.	1308
		1309
		1310
		1311
		1312
		1313
		1314
		1315
	Clara Meister, Ryan Cotterell, and Tim Vieira. 2020. If beam search is the answer, what was the question? In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 2173–2185, Online. Association for Computational Linguistics.	1316
		1317
		1318
		1319
		1320
		1321
	Haitao Mi, Zhiguo Wang, and Abe Ittycheriah. 2016. Supervised attentions for neural machine translation . In <i>Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing</i> , pages 2283–2288, Austin, Texas. Association for Computational Linguistics.	1322
		1323
		1324
		1325
		1326
		1327
	Paul Michel, Omer Levy, and Graham Neubig. 2019. Are sixteen heads really better than one? In <i>Advances in Neural Information Processing Systems</i> , volume 32. Curran Associates, Inc.	1328
		1329
		1330
		1331
	Grégoire Montavon, Sebastian Lapuschkin, Alexander Binder, Wojciech Samek, and Klaus-Robert Müller. 2017. Explaining nonlinear classification decisions with deep taylor decomposition. <i>Pattern Recognition</i> , 65:211–222.	1332
		1333
		1334
		1335
		1336

1337	Pooya Moradi, Nishant Kambhatla, and Anoop Sarkar.	Vikas Raunak, Arul Menezes, and Marcin Junczys-Dowmunt. 2021. The curious case of hallucinations in neural machine translation. <i>arXiv preprint arXiv:2104.06683</i> .	1394
1338	2020. Training with adversaries to improve faithfulness of attention in neural machine translation . In <i>Proceedings of the 1st Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 10th International Joint Conference on Natural Language Processing: Student Research Workshop</i> , pages 93–100, Suzhou, China. Association for Computational Linguistics.	1395	
1339		1396	
1340		1397	
1341			
1342		Abhilasha Ravichander, Yonatan Belinkov, and Eduard Hovy. 2021. Probing the probing paradigm: Does probing accuracy entail task relevance? In <i>Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume</i> , pages 3363–3377, Online. Association for Computational Linguistics.	1398
1343		1399	
1344		1400	
1345		1401	
1346	Rajiv Movva and Jason Zhao. 2020. Dissecting lottery ticket transformers: Structural and behavioral study of sparse neural machine translation . In <i>Proceedings of the Third BlackboxNLP Workshop on Analyzing and Interpreting Neural Networks for NLP</i> , pages 193–203, Online. Association for Computational Linguistics.	1402	
1347		1403	
1348		1404	
1349		1405	
1350		1406	
1351		1407	
1352		1408	
1353	Mathias Müller and Rico Sennrich. 2021. Understanding the properties of minimum bayes risk decoding in neural machine translation . In <i>Proceedings of the Joint Conference of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (ACL-IJCNLP 2021)</i> .	1409	
1354		1410	
1355		1411	
1356		1412	
1357		1413	
1358		1414	
1359		1415	
1360	Rafael Müller, Simon Kornblith, and Geoffrey E. Hinton. 2019. When does label smoothing help? In <i>NeurIPS</i> .	1416	
1361		1417	
1362		1418	
1363		1419	
1364	Kenton Murray and David Chiang. 2018. Correcting length bias in neural machine translation . In <i>Proceedings of the Third Conference on Machine Translation: Research Papers</i> , pages 212–223, Brussels, Belgium. Association for Computational Linguistics.	1420	
1365		1421	
1366		1422	
1367		1423	
1368			
1369		1424	
1370	Myle Ott, Michael Auli, David Grangier, and Marc’Aurelio Ranzato. 2018. Analyzing uncertainty in neural machine translation . In <i>Proceedings of the 35th International Conference on Machine Learning</i> , volume 80 of <i>Proceedings of Machine Learning Research</i> , pages 3956–3965. PMLR.	1425	
1371		1426	
1372		1427	
1373		1428	
1374		1429	
1375		1430	
1376	Alessandro Raganato, Yves Scherrer, and Jörg Tiedemann. 2020. Fixed encoder self-attention patterns in transformer-based machine translation . In <i>Findings of the Association for Computational Linguistics: EMNLP 2020</i> , pages 556–568, Online. Association for Computational Linguistics.	1431	
1377		1432	
1378			
1379		1433	
1380		1434	
1381	Marc’Aurelio Ranzato, Sumit Chopra, Michael Auli, and Wojciech Zaremba. 2015. Sequence level training with recurrent neural networks. <i>arXiv preprint arXiv:1511.06732</i> .	1435	
1382		1436	
1383		1437	
1384		1438	
1385		1439	
1386		1440	
1387	Vikas Raunak, Siddharth Dalmia, Vivek Gupta, and Florian Metze. 2020. On long-tailed phenomena in neural machine translation . In <i>Findings of the Association for Computational Linguistics: EMNLP 2020</i> , pages 3088–3095, Online. Association for Computational Linguistics.	1441	
1388		1442	
1389		1443	
1390		1444	
1391		1445	
1392		1446	
1393			
		Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. 2017. Grad-cam: Visual explanations from deep networks via gradient-based localization. In <i>Proceedings of the IEEE international conference on computer vision</i> , pages 618–626.	1447
			1448
		Jiajun Shen, Peng-Jen Chen, Matthew Le, Junxian He, Jiatao Gu, Myle Ott, Michael Auli, and Marc’Aurelio Ranzato. 2021. The source-target domain mismatch problem in machine translation . In	1449
			1450

1565		<i>Computational Linguistics</i> , pages 1198–1212, Florence, Italy. Association for Computational Linguistics.		
1566				
1567				
1568	Elena Voita, Rico Sennrich, and Ivan Titov. 2021. Language modeling, lexical translation, reordering: The training process of NMT through the lens of classical SMT. <i>CoRR</i> , abs/2109.01396.			
1569				
1570				
1571				
1572	Elena Voita, David Talbot, Fedor Moiseev, Rico Sennrich, and Ivan Titov. 2019b. Analyzing multi-head self-attention: Specialized heads do the heavy lifting, the rest can be pruned. <i>arXiv preprint arXiv:1905.09418</i> .			
1573				
1574				
1575				
1576				
1577	Chaghan Wang, Anirudh Jain, Danlu Chen, and Jitao Gu. 2019a. VizSeq: a visual analysis toolkit for text generation tasks. In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP): System Demonstrations</i> , pages 253–258, Hong Kong, China. Association for Computational Linguistics.			
1578				
1579				
1580				
1581				
1582				
1583				
1584				
1585				
1586	Chaojun Wang and Rico Sennrich. 2020. On exposure bias, hallucination and domain shift in neural machine translation. In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 3544–3552, Online. Association for Computational Linguistics.			
1587				
1588				
1589				
1590				
1591				
1592	Qiang Wang, Bei Li, Tong Xiao, Jingbo Zhu, Changliang Li, Derek F. Wong, and Lidia S. Chao. 2019b. Learning deep transformer models for machine translation. In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 1810–1822, Florence, Italy. Association for Computational Linguistics.			
1593				
1594				
1595				
1596				
1597				
1598				
1599	Shuo Wang, Zhaopeng Tu, Shuming Shi, and Yang Liu. 2020. On the inference calibration of neural machine translation. In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 3070–3079, Online. Association for Computational Linguistics.			
1600				
1601				
1602				
1603				
1604				
1605	Johnny Wei, Khiem Pham, Brendan O’Connor, and Brian Dillon. 2018. Evaluating grammaticality in seq2seq models with a broad coverage HPSG grammar: A case study on machine translation. In <i>Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP</i> , pages 298–305, Brussels, Belgium. Association for Computational Linguistics.			
1606				
1607				
1608				
1609				
1610				
1611				
1612				
1613	Sarah Wiegrefe and Yuval Pinter. 2019. Attention is not not explanation. In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 11–20, Hong Kong, China. Association for Computational Linguistics.			
1614				
1615				
1616				
1617				
1618				
1619				
			Sam Wiseman and Alexander M. Rush. 2016. Sequence-to-sequence learning as beam-search optimization. In <i>Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing</i> , pages 1296–1306.	1620
				1621
				1622
				1623
				1624
			Lijun Wu, Fei Tian, Tao Qin, Jianhuang Lai, and Tie-Yan Liu. 2018. A study of reinforcement learning for neural machine translation. In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 3612–3621, Brussels, Belgium. Association for Computational Linguistics.	1625
				1626
				1627
				1628
				1629
				1630
				1631
			Y. Wu, M. Schuster, Z. Chen, Quoc V. Le, Mohammad Norouzi, Wolfgang Macherey, M. Krikun, Yuan Cao, Q. Gao, Klaus Macherey, J. Klingner, Apurva Shah, M. Johnson, X. Liu, Lukasz Kaiser, Stephan Gouws, Y. Kato, Taku Kudo, H. Kazawa, K. Stevens, George Kurian, Nishant Patil, W. Wang, C. Young, J. Smith, Jason Riesa, Alex Rudnick, Oriol Vinyals, G. Corrado, Macduff Hughes, and J. Dean. 2016. Google’s neural machine translation system: Bridging the gap between human and machine translation. <i>ArXiv</i> , abs/1609.08144.	1632
				1633
				1634
				1635
				1636
				1637
				1638
				1639
				1640
				1641
				1642
			Ruibin Xiong, Yunchang Yang, Di He, Kai Zheng, Shuxin Zheng, Chen Xing, Huishuai Zhang, Yanyan Lan, Liwei Wang, and Tiejun Liu. 2020. On layer normalization in the transformer architecture. In <i>International Conference on Machine Learning</i> , pages 10524–10533. PMLR.	1643
				1644
				1645
				1646
				1647
				1648
			Weijia Xu, Shuming Ma, Dongdong Zhang, and Marine Carpuat. 2021. How does distilled data complexity impact the quality and confidence of non-autoregressive machine translation? In <i>Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021</i> , pages 4392–4400, Online. Association for Computational Linguistics.	1649
				1650
				1651
				1652
				1653
				1654
				1655
			Baosong Yang, Longyue Wang, Derek F. Wong, Lidia S. Chao, and Zhaopeng Tu. 2019. Assessing the ability of self-attention networks to learn word order. In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 3635–3644, Florence, Italy. Association for Computational Linguistics.	1656
				1657
				1658
				1659
				1660
				1661
				1662
			Yilin Yang, Liang Huang, and Mingbo Ma. 2018. Breaking the beam search curse: A study of (re-)scoring methods and stopping criteria for neural machine translation. <i>arXiv preprint arXiv:1808.09582</i> .	1663
				1664
				1665
				1666
			Kayo Yin, Patrick Fernandes, Danish Pruthi, Aditi Chaudhary, André F. T. Martins, and Graham Neubig. 2021. Do context-aware translation models pay the right attention? In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 788–801, Online. Association for Computational Linguistics.	1667
				1668
				1669
				1670
				1671
				1672
				1673
				1674
				1675

1676	Gal Yona, Amirata Ghorbani, and James Zou. 2021.	our surveyed papers. Please refer to Belinkov et al.	1728
1677	Who’s responsible? jointly quantifying the contri-	(2020b); Belinkov and Glass (2019) for a general	1729
1678	bution of the learning algorithm and data. In <i>Pro-</i>	introduction to interpretation methodologies.	1730
1679	<i>ceedings of the 2021 AAAI/ACM Conference on AI,</i>		
1680	<i>Ethics, and Society</i> , AIES ’21, page 1034–1041,		
1681	New York, NY, USA. Association for Computing		
1682	Machinery.		
1683	Thomas Zenkel, Joern Wuebker, and John DeNero.	B.1 Attribution	1731
1684	2020. End-to-end neural word alignment outper-	Attribution is one of the local explanation method-	1732
1685	forms GIZA++. In <i>Proceedings of the 58th Annual</i>	ologies for understanding and visualizing the deci-	1733
1686	<i>Meeting of the Association for Computational Lin-</i>	sion of predictive models, i.e., classifiers (Carvalho	1734
1687	<i>guistics</i> , pages 1605–1617, Online. Association for	et al., 2019). It relates every model prediction to	1735
1688	Computational Linguistics.	a subset of input features that might be the cause	1736
1689	Yuming Zhai, Gabriel Illouz, and Anne Vilnat. 2020.	of that prediction. A large number of attribution	1737
1690	Detecting non-literal translations by fine-tuning	methods are proposed recently in vision and learn-	1738
1691	cross-lingual pre-trained language models. In	ing community (Simonyan et al., 2013 ; Bach et al.,	1739
1692	<i>Proceedings of the 28th International Conference</i>	2015 ; Montavon et al., 2017 ; Selvaraju et al., 2017 ;	1740
1693	<i>on Computational Linguistics</i> , pages 5944–5956,	Sundararajan et al., 2017). In NMT, the predic-	1741
1694	Barcelona, Spain (Online). International Committee	tion \hat{y} could be seen as a sequence of classification	1742
1695	on Computational Linguistics.	steps. Given input x , predicted sequence \hat{y} , and the	1743
1696	Wen Zhang, Yang Feng, Fandong Meng, Di You, and	NMT model \mathcal{M} , an attribution method is defined	1744
1697	Qun Liu. 2019. Bridging the gap between training	as an algorithmic process $\mathcal{A}(x, \hat{y}, \mathcal{M})$, it outputs	1745
1698	and inference for neural machine translation. <i>arXiv</i>	<i>relevant scores</i> over every token of x and $y_{<t}$ for	1746
1699	<i>preprint arXiv:1906.02448.</i>	each \hat{y}_t . Based on relevant scores, we can at least	1747
1700	Zhengli Zhao, Dheeru Dua, and Sameer Singh. 2017.	qualitatively know what \hat{y}_t is probably aligned to.	1748
1701	Generating natural adversarial examples. <i>CoRR</i> ,	Model-specific Methods Model-specific attribu-	1749
1702	abs/1710.11342.	tions can have access to the inner computation of	1750
1703	Chunting Zhou, Jiatao Gu, Mona T. Diab, P. Guzmán,	the NMT model. Gradient-based attribution uses	1751
1704	Luke Zettlemoyer, and Marjan Ghazvininejad.	the activation of \hat{y}_t for backward computation. It	1752
1705	2020a. Detecting hallucinated content in condi-	then uses the gradients on each embedding vector	1753
1706	tional neural sequence generation. <i>ArXiv</i> ,	of every token in x and $\hat{y}_{<t}$ to compute its rele-	1754
1707	abs/2011.02593.	vant score regarding \hat{y}_t (Ding et al., 2019). Layer-	1755
1708	Chunting Zhou, Jiatao Gu, and Graham Neubig.	wise Relevance Propagation (LRP) uses activation	1756
1709	2020b. Understanding knowledge distillation in	vectors and model weights to compute neuron rele-	1757
1710	non-autoregressive machine translation. In <i>8th Inter-</i>	vance, and then back-propagates the relevance back	1758
1711	<i>national Conference on Learning Representations,</i>	to the input tokens (Ding et al., 2017).	1759
1712	<i>ICLR 2020, Addis Ababa, Ethiopia, April 26-30,</i>		
1713	<i>2020.</i> OpenReview.net.	Model-agnostic Methods Model-agnostic attri-	1760
1714	Conghui Zhu, Guanlin Li, Lemao Liu, Tiejun Zhao,	butions regard the NMT model as a black-box.	1761
1715	and Shuming Shi. 2020. Understanding learning dy-	These methods usually calculate the relevant scores	1762
1716	namics for neural machine translation.	through manipulating model inputs (Alvarez-Melis	1763
1717	A Mindmap of FAQs	and Jaakkola, 2017 ; Li et al., 2019b). For example,	1764
1718	Fig. 1 demonstrates a screenshot of the mindmap on	prediction difference (Li et al., 2019b) chooses a	1765
1719	our website (https://nmtology.github.	particular feature (token $x_{t'}$), and observe the prob-	1766
1720	io/). Visitors can zoom in or zoom out the tree by	ability difference of \hat{y}_t before and after removing	1767
1721	the clicking inner nodes. And by clicking a specific	that feature, the larger the probability is, the more	1768
1722	question, you will be guided to a separate webpage	relevant between \hat{y}_t and the removed one.	1769
1723	that hosts the related papers under that question.		
1724	B Methodology	B.2 Probing	1770
1725	This section gives a focused introduction to several	Probing is a method for investigating how much	1771
1726	commonly used methodologies for understanding	a component of the NMT model captures certain	1772
1727	the NMT framework, which are commonly used in	kind of knowledge. The main technique for probing	1773
		is to train a classifier g which maps an intermedi-	1774
		ate representation $f(x)$ of the input x to certain	1775
		property of interest z (Alain and Bengio, 2016 ;	1776

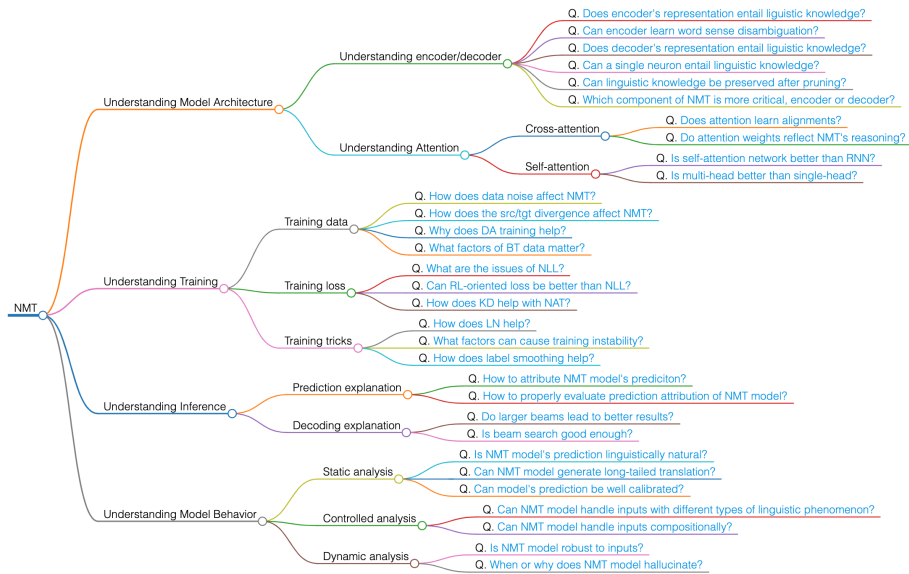


Figure 1: A screenshot of the mindmap of FAQs on our website <https://nmtology.github.io/>.

Conneau et al., 2018). This network component $f(\cdot)$ can be word embedding, sentence embedding, hidden state, attention weight, etc. The property z can be various linguistic features, such as part-of-speech tags, morphological information, or more complicated syntactic or semantic features. Then, the accuracy of $g(f(x))$ can reveal the quality of representations $f(x)$ with respect to the property z , so that different model components can be compared to each other. To show this accuracy is non-trivial, it can be compared to feeding random inputs into the classifier $g(\cdot)$. Meanwhile, comparing with state-of-the-art on that task can inform us how much is missing from the representation.

B.3 Others

In addition to attribution and probing which are most commonly used, several other methodologies are used in specific analysis scenarios (Blinkov and Glass, 2019). *i) Visualization* is always used accompanied with attribution to show the relationship between predicted and input tokens; it is also used to visualize clustering effects of learned representations through dimension reduction techniques (Alvarez-Melis and Jaakkola, 2017; Ding et al., 2017). *ii) Challenge set* is always used to investigate certain desirable characteristic of the model through data test suite construction (Isabelle et al., 2017). *i) Model extraction* is to extract use knowledge distillation to learn a transparent or interpretable surrogate NMT model (e.g. rules, syntactic trees) from the original one (Bastani et al.,

2017; Sushil et al., 2018). Besides, several works also build toolkits for visualization and model debugging (Strobel et al., 2018; Wang et al., 2019a; Graliński et al., 2019; Gauthier et al., 2020).

1777
1778
1779
1780
1781
1782
1783
1784
1785
1786
1787
1788
1789
1790
1808
1809
1810
1811