Al in GILT: A Decade of Developments (2014–2024)

GILT Research Centre (https://giltrc.org) The Hang Seng University of Hong Kong 2024

The present work is supported by a grant from the Research Grants Council of the Hong Kong Special Administrative Region, China (Project Reference No.: UGC/ IDS(R)14/23).



Introduction

The GILT acronym represents the combination of Globalisation, Internationalisation, Localisation, and Translation practices that enable content and products to reach global audiences through multilanguage adaptation. The language industry created the term GILT which defines "globalization" as business strategies for international reach while "internationalization" (i18n) describes product design for localization ease and "localization" (I10n) means adapting content technically and culturally for different locales and "translation" represents converting text between languages. According to Cadieux and Esselink (2002), translation stands as the most familiar term among these while the industry has been continually refining globalization, internationalization, and localization definitions and practices [1].

In the past decade, artificial intelligence (AI) – particularly advances in natural language processing – have dramatically transformed GILT. This academic monograph provides a chronological review of key English-language publications from the last ten years (2014–2024) that have shaped AI in GILT and GILT practices themselves. We cover breakthroughs in machine translation (MT), which lies at the heart of localization, including the shift from statistical approaches to neural networks and the emergence of transformative architectures (e.g. the Transformer [2]). We also examine research on integrating MT into localization workflows, human-computer interaction in translation tools, and the rise of large language models (LLMs) like ChatGPT that have begun to impact global communication. Foundational AI papers (such as Attention Is AII You Need in 2017 [2]) are included for their outsized influence on language technology.

The review is organized by year, with each section summarizing a significant publication. Each subsection provides the paper's contributions, methodology, and its significance for GILT in an academic context. Full bibliographical details are given in APA format for each work. Through this chronological narrative, we trace how AI innovations have advanced translation quality, expanded language coverage, improved tools for translators, and influenced the strategies and capabilities in globalization and localization.

2014: Neural Machine Translation and Interactive Translation Emerge

Cho et al. (2014) – Learning Phrase Representations using RNN Encoder–Decoder for SMT

In 2014, Cho and colleagues introduced one of the first neural network models for machine translation, pioneering the RNN Encoder–Decoder architecture [3]. Their model consists of two recurrent neural networks: an encoder that reads a source sentence and compresses it into a vector, and a decoder that generates a translation from that vector. This joint encoder–decoder is trained end-to-end to maximize the conditional probability of the target (translation) given the source sentence [3]. Notably, the model learns continuous vector representations of phrases that capture syntactic and semantic properties, enabling it to handle phrases more flexibly than phrase-based statistical MT. Incorporating the RNN Encoder–Decoder as a feature in a standard statistical MT system yielded improved translation performance [3]. Qualitatively, the authors showed that the neural model learned meaningful representations of linguistic phrases [3]. Significance: This paper demonstrated the feasibility of purely neural translation components and introduced the concept of jointly learning to encode and decode sequences, laying groundwork for subsequent neural machine translation (NMT) research.

Sutskever et al. (2014) – Sequence to Sequence Learning with Neural Networks

Sutskever et al. built on emerging encoder–decoder ideas by creating a general end-to-end sequenceto-sequence (seq2seq) learning framework using deep LSTM networks. Their seq2seq model employs one LSTM to encode an input sequence (e.g., a sentence) into a fixed-length vector, and a second LSTM to decode that vector into an output sequence [4]. Importantly, this approach makes minimal assumptions about the sequence structure [4], allowing the model to learn directly from data. In their experiments on English-to-French translation, a 4-layer LSTM seq2seq model achieved impressive results, especially when using a technique called training with reversed source sequences to ease optimization. Contributions: This work provided a simple yet powerful recipe for training neural networks to map sequences to sequences, and it showed for the first time that an RNN with sufficient capacity could learn to translate entire sentences. Significance for GILT: The seq2seq paradigm became the foundation of NMT. By late 2014, neural models (Cho et al. [3] and Sutskever et al. [4]) had established a new direction for MT, promising more fluent and generalizable translations than the thendominant phrase-based systems.

Green et al. (2014) – Human Effort and Machine Learnability in Computer-Aided Translation

While neural MT was emerging, another 2014 study by Green et al. examined how AI could assist human translators via interactive translation. They presented a new translator workstation that allowed for two modes of assistance: traditional post-editing of machine output, and a novel interactive MT mode where the system offers suggestions as the translator types [5]. They quantitatively evaluated professional translators using both modes on English-French and English-German tasks. The results showed that post-editing yielded faster throughput, but interactive translation produced slightly higher quality outputs [5]. Moreover, they introduced an adaptive MT component: the underlying MT engine was incrementally retrained ("re-tuned") on the translator's corrections in real time. Interestingly, retraining on the interactive mode's data led to significantly larger quality gains (measured by HTER reduction) than retraining on post-edited data [5]. Contributions: This work provided the first holistic comparison of post-editing vs. interactive translation and proposed methods for MT systems to learn from human feedback on the fly. Significance for GILT: It highlighted human-centric design in translation tools and showed that with the right interface and adaptive algorithms, MT can effectively augment human translators, foreshadowing later developments in translator productivity tools and interactive neural MT.

2015: Attention Mechanisms

Bahdanau et al. (2015) – Neural Machine Translation by Jointly Learning to Align and Translate

Bahdanau et al. tackled a key limitation of early seq2seq models – the fixed-length bottleneck – by introducing an attention mechanism into neural MT [6]. In their approach, the decoder doesn't attempt to compress all information into a single vector; instead, at each translation step it automatically learns to align (focus attention on) the parts of the source sentence most relevant to predicting the next target word [6]. This soft alignment process effectively lets the model "look back" at the source sentence as needed, which greatly improves translation of long sentences. Using the new attention-based model, they achieved translation quality on an English–French task comparable to the state-of-the-art phrase-based system of the time [6]. Just as importantly, the learned attention weights corresponded closely to human-like word alignments [6], providing an interpretable insight into how the model translated. Significance: This paper's attention mechanism revolutionized NMT. It enabled models to handle long inputs and significantly improved translation accuracy. Nearly all subsequent NMT and NLP architectures – including the Transformer in 2017 [2] – built upon the idea of attention to better capture relationships in sequences. For localization, Bahdanau's attention model meant MT systems could more reliably translate complex sentences without losing context, a critical improvement for quality and fidelity.

Zaretskaya et al. (2015) – Integration of Machine Translation in CAT Tools: State of the Art, Evaluation, and User Attitudes

Zaretskaya and colleagues shifted focus to the practical integration of MT into Computer-Assisted Translation (CAT) tools used by human translators. This study surveyed and evaluated the state-of-theart in 2015 for incorporating MT suggestions into translation memory (TM) workflows. It examined major translation environments and how they offered MT proposals to translators, and it gathered feedback from professional translators on using MT within their CAT interfaces. The findings indicated that when MT is integrated seamlessly (for example, by displaying MT suggestions alongside TM matches), translators can benefit from increased productivity and maintain accuracy [7]. However, user attitudes were mixed and depended on MT quality and the effort needed to post-edit. Many translators were cautious – some appreciated MT as a helpful reference especially for low-fuzzy matches, while others reported frustration when MT output was poor or when it disrupted their typical workflow. Significance: This work provided an early snapshot of how MT was being adopted on the translator's desktop. It highlighted usability challenges and the importance of user-centered design in MT integrations. For the GILT community, it underscored that even as MT quality was improving, human translators' acceptance of MT depended on tool ergonomics and trust in MT output.

2016: Handling Vocabulary, Data Augmentation, and Google's Leap to Neural MT

Sennrich et al. (2016a) – Neural Machine Translation of Rare Words with Subword Units

As neural MT took off, one practical problem became apparent: fixed vocabularies. Early NMT systems could struggle with out-of-vocabulary words (e.g., proper names, rare terms) because they treated words as indivisible tokens. Sennrich et al. addressed this by using subword units learned through a data-driven compression algorithm (byte-pair encoding, or BPE). They showed that by splitting words into smaller units (like pieces of words or syllables), an NMT model can achieve open-vocabulary translation, effectively handling rare or unknown words by constructing them from subword components [8]. In experiments on English→German/Russian, the BPE-based model outperformed a baseline that backed off to dictionary translations for unknown words, improving BLEU by ~1.1−1.3 points [8]. Significance: This paper provided a simple and robust solution for one of NMT's early weaknesses. Subword modeling (with BPE or similar methods) quickly became a standard in MT and NLP, enabling models to cope with the rich morphology of many languages and to generate terminology not seen in training. For the localization industry, this meant neural MT could cover product names, technical terms, and inflections more reliably, an important step for practical deployment.

Sennrich et al. (2016b) – Improving Neural Machine Translation Models with Monolingual Data

A second breakthrough by Sennrich and colleagues in 2016 was demonstrating how to leverage monolingual target-language data to improve NMT – a technique now known as back-translation. They generated synthetic parallel training data by translating target-language sentences back into the source language using a baseline MT system [9]. These synthetic source–target pairs were then added to the training set. Using this method, they achieved substantial gains: e.g., +2.8 to +3.7 BLEU on English→German and +2.1 to +3.4 BLEU on Turkish→English, reaching new state-of-the-art results on those tasks [9]. They also showed that fine-tuning on in-domain monolingual and parallel data further improved performance. Significance: This work was pivotal in showing that NMT could utilize the abundant monolingual texts available, alleviating parallel data scarcity. Back-translation became a standard practice in MT model training, crucial for improving fluency and handling domain-specific language. In a GILT context, this meant MT systems could be adapted better to target-language style and terminology (using customer-specific monolingual texts, for example), enhancing the quality of localized translations.



Wu et al. (2016) – Google's Neural Machine Translation System: Bridging the Gap between Human and Machine Translation

In late 2016, Google announced its full switch from phrase-based to neural MT with the Google Neural Machine Translation (GNMT) system [10]. Wu et al. detailed this production-scale NMT: an 8-layer LSTM with residual connections and an attention mechanism, enhanced with engineering techniques for speed (e.g., low-precision arithmetic) and a method of splitting words into wordpieces (similar to BPE) to handle rare words [10]. The GNMT system obtained translations on WMT benchmarks on par with the best prior systems [10]. In a noteworthy human evaluation, Google reported GNMT reduced translation errors by an average of 60% compared to their existing phrase-based system on a set of simple test sentences [10]. This was a dramatic quality jump recognized as bridging much of the gap to human translation. Significance: Google's deployment of GNMT in production was a milestone for the entire localization industry – it signaled that NMT was not just a research novelty but ready for real-world, global-scale translation. GNMT's architecture and training innovations (such as massive parallel computing and wordpiece models) influenced many subsequent systems. For content globalization, Google's move meant billions of users suddenly experienced better machine translations (e.g., in Google Translate), raising expectations and demand for high-quality MT in many languages.

2017: The Transformer Revolution, Multilingual MT, and Human-Focused Studies

Vaswani et al. (2017) – Attention Is All You Need

Vaswani et al. introduced the Transformer, a novel neural architecture that did away with recurrence and convolutions entirely, relying only on self-attention mechanisms for sequence modeling [2]. The Transformer model processes all words in a sentence in parallel, using multi-head self-attention to learn dependencies between words regardless of their distance, and a feed-forward network at each position. This design enabled much greater training parallelization compared to RNNs. On translation tasks (English-German and English-French), Transformers not only trained faster but also achieved a new state-of-the-art: for English-German, a single Transformer model scored 28.4 BLEU, over 2 BLEU points better than the previous best ensemble [2]. On English-French, it established a new single-model record of 41.8 BLEU after training for just 3.5 days on eight GPUs, a small fraction of the training cost of earlier models [2]. Significance: The Transformer guickly became the dominant architecture in MT and NLP, thanks to its efficiency and effectiveness. For the GILT domain, Transformers meant MT systems could be scaled to very large datasets and models, yielding higher quality translations across many languages. This architecture also paved the way for the later development of large pretrained language models and multilingual Transformers that have greatly impacted global communication technology.

Johnson et al. (2017) – Google's Multilingual Neural Machine Translation System: Enabling Zero-Shot Translation

Johnson et al. (a Google team) extended NMT to a multilingual setting. They designed a single NMT model to translate between multiple languages by prepending a special token indicating the target language to each source sentence [11]. Using a shared wordpiece vocabulary, their single model was jointly trained on data for many language pairs. The results were remarkable: one multilingual model could match or exceed the quality of individual bilingual models for high-resource pairs like English– French and English–German [11]. Moreover, the model demonstrated an emergent ability to perform zero-shot translation – translating between language pairs it had never seen explicitly during training (e.g., translating directly from Japanese to Korean after being trained on Japanese–English and English–Korean) [11]. This suggested the model learned some universal interlingua representations. Significance: This was a milestone in massively multilingual AI. It showed that one model can handle dozens of languages, which is extremely valuable for globalization: instead of maintaining separate MT systems for each language pair, an organization could have one unified model. The zero-shot capability was particularly exciting for covering low-resource language combinations. In practice, Google deployed such multilingual models, leading to improvements in translation quality for many lesser-served languages and simplifying their MT infrastructure.

Koehn & Knowles (2017) – Six Challenges for Neural Machine Translation

Mid-2017, Koehn and Knowles provided a sober assessment of NMT's weaknesses relative to the established phrase-based MT. They identified six key challenges where NMT struggled or required attention: (1) handling out-of-domain data, (2) training with limited amounts of data, (3) translating rare words, (4) very long sentences, (5) word alignment transparency, and (6) effective beam search decoding [12]. For example, they showed NMT quality drops sharply when input sentences are much longer than those seen in training, and low-resource language pairs remained a major pain point. They also noted that while NMT's attention mechanism provides some alignment information, it is not as easily interpretable or controllable as alignment in phrase-based MT. Significance: This analysis was influential in guiding research priorities. It made clear that despite NMT's rapid advances, there was still "no free lunch" and plenty of work needed to adapt NMT to real-world conditions (e.g., low-resource languages, robust decoding). Many subsequent improvements – such as better training regularization, unsupervised and transfer learning for low-resource MT, and coverage models for attention – can be traced to the challenges outlined here. For the GILT community, Koehn & Knowles's paper tempered hype with practical insight, reminding practitioners that new NMT systems had to be evaluated carefully, especially for languages or scenarios with sparse data.

Moorkens & O'Brien (2017) – Assessing User Interface Needs of Post-Editors of Machine Translation

Moorkens and O'Brien turned the focus to the human translators who post-edit MT output, investigating what interface features they need for an effective and ergonomic experience [13]. Through interviews and observations of professional translators, they identified pain points in traditional CAT tool interfaces when post-editing, such as lack of fluency in the MT suggestions display, inadequate support for previewing the source text context, and insufficient customizability of the editing environment. Post-editors expressed desire for interface improvements like better visual cues for MT vs. TM suggestions, one-click reflows of suggestions after corrections, and predictive text that could speed up editing. This study also underscored the cognitive load differences between translating from scratch and post-editing, suggesting UI designs that reduce friction (for instance, by minimizing cursor movements or mode switching) could make post-editing more efficient and less "irritating" [13]. Significance: This work is an example of aligning technology with human factors. As MT became ubiquitous in localization workflows by 2017, understanding the needs of translators who must work with MT output was critical. The recommendations from Moorkens & O'Brien influenced the development of more user-friendly post-editing environments (such as dynamic quality estimation cues, improved shortcut workflows, etc.). It reinforced the idea that improving raw MT quality isn't the only goal - improving the human-machine interaction in translation tasks is equally important for overall efficiency and translator satisfaction.

2018: Milestones in MT Quality and New Frontiers (Unsupervised MT, Pretrained Language Models)

Hassan et al. (2018) – Achieving Human Parity on Automatic Chinese–English News Translation

In 2018, a team at Microsoft made headlines by declaring that their latest Chinese–English NMT system had achieved human parity on a popular news translation test set [14]. Using an ensemble of Transformer models plus extensive input pre-processing and model optimization, they reported that professional human evaluators could not distinguish the system's translations from human translations on the newstest2017 Chinese→English news dataset [14]. They defined "human parity" in a specific sense – that the system's translation quality on that test, under those conditions, was statistically indistinguishable from that of human translators. This announcement was a major publicity moment, demonstrating how far NMT had come in just a few years. Significance: If taken at face value, human parity suggests MT has reached a quality level suitable for certain professional contexts. For the GILT community, such a result implied that for at least some content types (e.g., structured news sentences in a high-resource language pair), MT output might be on par with human translation, potentially shifting workflows towards heavier use of raw MT. However, this claim also sparked discussion and follow-up research to verify and nuance the "human parity" assertion, illustrating the importance of rigorous evaluation beyond headline metrics.

Toral et al. (2018) – Attaining the Unattainable? Reassessing Claims of Human Parity in Neural MT

Toral et al. responded directly to Microsoft's human parity claim by conducting a careful analysis. They pointed out that the evaluation of MT vs. human had to consider factors like the origin of the test sentences (original language vs. translations) and the expertise of evaluators [15]. By re-evaluating the same Chinese→English outputs with professional translators as judges and focusing only on sentences originally written in Chinese (to avoid source-language bias), Toral and colleagues found that the human translators still outperformed the MT system [15]. In other words, when controlling for these variables, human parity was not actually achieved [15]. They also uncovered specific quality issues in the human reference translations and in the MT outputs that were overlooked in the initial study. Significance: This work injected nuance and scientific rigor into the evaluation of "human parity." It underscored that declaring MT equal to humans depends heavily on how evaluation is done. For practitioners, the lesson is that even very high-performing MT should be tested in realistic settings (with domain-expert reviewers, and using original texts) before drawing conclusions. The Toral et al. study helped establish more robust human evaluation methodologies for MT going forward, including at subsequent WMT evaluation campaigns.

Lample et al. (2018) – Unsupervised Machine Translation Using Monolingual Corpora Only

A major stride for low-resource MT came from Lample et al., who demonstrated unsupervised MT training a translation model with zero parallel sentences [16]. They leveraged only monolingual data in two languages, using two key ideas: (1) learning strong bilingual word embeddings (shared latent representations) to initialize a shared encoder-decoder, and (2) applying a cycle of back-translation in training (generating synthetic translations in both directions to create pseudo-parallel data). Evaluating on French-English (a language pair with abundant data to benchmark against), their unsupervised system achieved BLEU scores of 15.1 on WMT English \rightarrow French, remarkable given it saw no parallel training data [16]. While this was below supervised MT, it was far better than word-by-word translation. On a simpler dataset (Multi30k images captions, English-French), unsupervised MT reached 32.8 BLEU [16], approaching the quality of supervised models. Significance: This was a breakthrough for translation of low-resource languages or domains where parallel corpora are scarce. It opened a research line where many improved unsupervised and semi-supervised MT approaches followed. For GILT, unsupervised MT offers hope for extending translation technology to languages and niches with little translated content available (for example, it suggests a method to produce MT systems for languages of emerging markets by leveraging monolingual web data in those languages and a lingua franca). By late 2018, the community began to see that even without explicit bilingual data, AI could learn to translate - a powerful concept for global language inclusion.

2019: Advances in Language Modeling and Human–Machine Interaction

Devlin et al. (2019) – BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding

Although not an MT paper per se, BERT (Bidirectional Encoder Representations from Transformers) was a seminal 2018–2019 development in NLP that greatly influenced translation and localization technology [17]. Devlin et al. pre-trained a deep Transformer encoder on massive amounts of English text using a masked language modeling objective (predicting hidden words) and a next sentence prediction objective. The result was a contextual language model that could be fine-tuned to achieve state-of-the-art performance on a wide range of language tasks. BERT and its multilingual variant (M-BERT, trained on 104 languages) provided powerful cross-lingual representations that could be used for tasks like machine translation quality estimation, bilingual text classification, and terminology extraction. For example, multilingual BERT showed surprising ability to align representations across languages, enabling zero-shot cross-lingual transfer on tasks like question answering. Significance: BERT kicked off the era of large pre-trained language models. In the GILT context, this meant tools like translation memory matchers, content analyzers, and even MT systems could leverage pre-trained knowledge for better accuracy. Shortly after BERT, researchers integrated similar pre-training ideas into MT (via sequence-to-sequence pre-training like in Facebook's MASS and Google's mT5). BERT also improved many sub-tasks in localization workflows - from grammar checking to semantic search in multilingual content - making it a cornerstone of AI innovations in the translation industry at the close of the decade.

Radford et al. (2019) – Language Models are Unsupervised Multitask Learners

OpenAl's 2019 report on GPT-2 marked a leap in what generative language models could do [18]. GPT-2 was a large Transformer-based language model (with up to 1.5 billion parameters) trained on a huge corpus of web text in an unsupervised manner. The surprising finding was that this single model, without task-specific training, could generate coherent paragraphs of text and perform rudimentary reading comprehension, translation, and summarization purely from being prompted (this is the origin of the phrase "unsupervised multitask learners") [18]. For example, given a prompt in French, GPT-2 would continue in French; given a prompt and instructions to translate, it could produce a simple translation. While GPT-2's translation ability was far from state-of-the-art MT, the experiment showed that large language models acquire some translation competency as a byproduct of unsupervised training on multilingual web data. Significance: GPT-2 (and its successor GPT-3 in 2020) ushered in the era of Large Language Models (LLMs). For the GILT field, this development hinted at a future where extremely large general models might contribute to translation tasks, either by generating translations directly or by assisting MT systems through synthetic data generation and advanced language understanding. It also foreshadowed tools like ChatGPT, which by 2023 would be capable of fairly robust translation among its many other tasks – all without explicit bilingual training data.

Daems & Macken (2019) – Interactive Adaptive SMT versus Interactive Adaptive NMT: A User Experience Evaluation

As neural MT began replacing phrase-based MT in tools, Daems and Macken evaluated how translators fared with an interactive translation prediction tool using NMT vs. the older SMT [19]. In their study, professional translators used a CAT tool (LILT) in two configurations: one backed by an adaptive SMT engine and one by an adaptive NMT engine. Both engines could learn from the user's edits in real time (adapting to domain/style). The user experience metrics - including translation quality of final output, time taken, keystrokes, and qualitative satisfaction - were compared. The findings showed that the interactive NMT system led to higher post-editing productivity and was preferred by translators on average, primarily due to the more fluent suggestions it provided. The NMTbased suggestions required fewer corrections, improving the ergonomics of the interaction. However, the study also noted that when NMT made errors, they could be more subtle (and thus sometimes missed by the translator) compared to SMT's often overt mistakes. Significance: This was one of the first studies to directly examine translator interactions with NMT in a real work scenario. It provided evidence that NMT can yield not just better automatic scores, but also tangible improvements in human translator efficiency and experience when integrated into tools. For industry, this reinforced the push to upgrade CAT tools and translation workflows to neural backends, and it highlighted the need for training translators to be aware of NMT's strengths and weaknesses in an interactive setting.

2020–2021: Scaling Up – GPT–3 and Massively Multilingual Models

Brown et al. (2020) – Language Models are Few-Shot Learners

OpenAI's GPT-3 dramatically expanded on the scale of GPT-2, with 175 billion parameters, and showed an unprecedented ability to perform NLP tasks in a "fewshot" manner - providing only a few examples in the prompt without any fine-tuning [20]. For translation specifically, GPT-3 demonstrated that with appropriate prompting (e.g. giving a couple of example translations), it could translate between many languages to a reasonable level. In their paper, Brown et al. reported GPT-3 achieved respectable BLEU scores on translation benchmarks in a zero/few-shot setting, often approaching the performance of supervised NMT for high-resource languages. However, it still fell short of state-of-the-art dedicated MT systems, especially for more complex or low-resource language pairs. Significance: GPT-3 reinforced the trend of large general models encroaching on traditionally separate tasks like translation. It suggested a future workflow where a single AI model might handle diverse tasks (translating, summarizing, answering questions) on the fly. For GILT, GPT-3's emergence meant that the line between "translation engine" and "language engine" began to blur. While specialized MT models remained superior in 2020, the idea that a general AI could translate reasonably well on demand was revolutionary. This also influenced industry thinking – by late 2020, companies started to explore using large pre-trained models to augment translation pipelines (for example, using GPT-3 to generate alternative phrasings or to translate when MT systems lacked certain language pairs).



Fan et al. (2021) – Beyond English-Centric Multilingual Machine Translation

(Meta AI's work initially released in 2020) Fan et al. pushed multilingual MT to a new level by training a single massive model (called M2M-100) that can directly translate between 100 languages without pivoting through English [21]. Crucially, they collected and curated a huge multilingual corpus with many language pairs that do not include English (to avoid an English-centric bias). The resulting Transformer model, with 15 billion parameters, learned to translate any of the 100 languages to any other. On many low-resource language pairs, it achieved large guality improvements over English-pivot baselines, and even for richer languages it often outperformed pivoting. For example, on Chinese→French, the direct M2M model surpassed a system that translated Chinese→English→French. Significance: This project demonstrated the feasibility of a true "universal translator" model covering dozens of languages in a single network. For global localization efforts, such a model has clear appeal: it simplifies deployment (one model instead of many) and improves quality for non-English content translation. It also emphasizes inclusivity - languages that were previously sidelined in MT research received more attention under this approach (the motto "no language left behind," which Meta would adopt in a subsequent 2022 project, encapsulates this spirit). By 2020, the era of English-centric translation technology was beginning to give way to a more multilingual-focused paradigm, vital for regions where English is not the hub language.

2022: Democratizing AI for Translation – Instruction–Tuned Models and 200– Language MT

Ouyang et al. (2022) – Training Language Models to Follow Instructions with Human Feedback

Ouyang et al. described the techniques behind InstructGPT, a model that would form the backbone of ChatGPT [22]. The key contribution was Reinforcement Learning from Human Feedback (RLHF): after pre-training a large language model, they fine-tuned it to better follow user instructions by collecting demonstrations and preference comparisons from human annotators [22]. Although the paper's focus was not specifically on translation, the resulting model (an instruction-following GPT-3 variant) showed greatly improved usability in interactive settings. By late 2022, OpenAI deployed this as ChatGPT, which quickly proved capable of producing useful translations among its many capabilities. Significance: The instruction-tuning paradigm made large language models much more aligned to what users ask, meaning a system like ChatGPT can seamlessly translate when prompted in plain language (e.g., "Please translate this paragraph into Spanish"). For the translation/localization industry, ChatGPT's emergence in 2022 was a watershed moment - it introduced a conversational, on-demand translation tool that could handle many languages fairly well, without being explicitly an MT system. This spurred discussions on how such AI might be used: as a translator's assistant (generating draft translations or providing alternatives), or even for end-user facing translation in certain scenarios. Ouyang et al.'s work thus bridged advances in language understanding with practical translation use cases by making AI more user-friendly and controllable.

NLLB Team (2022) – No Language Left Behind: Scaling Human-Centered Machine Translation

In 2022, Meta AI announced the No Language Left Behind (NLLB) project, which produced an MT model covering an unprecedented 200 languages [23]. The NLLB paper details how the team tackled data collection and guality for many low-resource languages (through a combination of web mining, manual data curation, and creating a dedicated evaluation benchmark called FLORES-200). The resulting model, built with a 54-billion-parameter Transformer with sparse gating (Mixture-of-Experts), delivered high-quality translations even for languages with little prior MT support [23]. On average, NLLB improved BLEU by 44% over previous state-of-the-art on a 200-language test set [23]. Critically, they also evaluated the performance of over 40,000 translation directions using human translated test sets, and combined human evaluation with a novel toxicity benchmark covering all languages to assess translation safety. Significance: NLLB represents a major step toward Al inclusion - extending good translation to communities that were previously left behind in the digital language divide. For globalization, this means the potential to localize content into dozens of African, Asian, and minority languages that companies had rarely targeted due to lack of technology. The work also set new standards in evaluation: by creating human-translated test sets for 200 languages, it provided a way to actually measure progress for those languages. NLLB's open release of models and data in 2022 enabled researchers and practitioners worldwide to build upon their 200-language model, directly contributing to more inclusive localization efforts.

Martikainen (2022) – Ghosts in the Machine: Can Adaptive MT Help Reclaim a Place for the Human in the Loop?

Hanna Martikainen's 2022 study reflected on the integration of adaptive MT in professional translation training and work [24]. By having translation students use an interactive, adaptive MT tool (LILT) over a semester, she gathered insights on how this technology affects their perception of the translation process. Students acknowledged that adaptive MT (which learns from their edits in real time) can make them feel "more in control" compared to static MT, because the system's suggestions improve as they translate [24]. However, the study also noted that ultimate usability still depended more on core MT quality and the overall CAT tool ergonomics than on adaptivity alone [24]. Martikainen contextualized her experiment with industry trends, citing that early claims of MT reaching or surpassing human translators (e.g., Google in 2016 [10], Microsoft in 2018 [14]) were often overgeneralized [24]. She argued that the real promise of adaptive MT is in fostering a true human-AI partnership, but that requires tools specifically designed to empower the human translator, not sideline them. Significance: This work is representative of the translation community's response to AI advances - rather than passively accepting "human parity" narratives, translators and researchers are actively exploring how to shape technology so that human linguists remain central. By 2022, the notion of a "human in the loop" was paramount: Martikainen's findings support that adaptive MT can be one way to keep translators engaged and supported by AI, rather than feeling replaced by it. Such research influences how translation companies implement MT (e.g., offering training and adaptive systems rather than one-size-fits-all MT) and how translators are taught to work effectively with AI assistance.

2023–2024: Evaluating AI vs. Human Translation in the Era of LLMs

Yan et al. (2023) – GPT-4 vs. Human Translators: A Comprehensive Evaluation

As large language models like GPT-4 reached unprecedented levels, researchers began rigorously comparing them to professional translators. Yan et al. conducted a comprehensive evaluation where GPT-4's translations were pitted against those of human translators across multiple language pairs and topical domains [25]. They used carefully designed human evaluation: multiple rounds of finegrained error annotation to assess accuracy and fluency. The study found that GPT-4's quality was comparable to that of junior professional translators, though it still lagged behind experienced senior translators in some aspects of accuracy [25]. Specifically, GPT-4 tended to make errors like overly literal translations or subtle omissions that skilled humans would avoid, and its performance declined on resource-poor language directions [25]. Nonetheless, the fact that GPT-4 could achieve near-human performance in certain settings (and even surpass less experienced humans) is remarkable. Significance: By 2023, it became clear that for certain content and language pairs, LLMs like GPT-4 can deliver translations that are nearly publishable with minimal editing. This raises both opportunities and challenges: translation companies might use GPT-4 to boost productivity (assigning human translators more of a review/editorial role), but it also demands new evaluation standards and quality control processes, since LLMs may fail unpredictably. Yan et al.'s work exemplifies how the community adapted by developing more granular evaluation methodologies to pinpoint where AI translations diverge from human expectations. It underscores that even if "AI translators" are not perfect, the gap has narrowed to the point that integration and careful oversight of such tools can yield significant efficiency gains in localization workflows.

OpenAl (2023) – GPT-4 Technical Report

OpenAI's March 2023 technical report on GPT-4 provided an overview of this model's broad capabilities and limitations [26]. Of note for translation, GPT-4 demonstrated substantially improved multilingual abilities over its predecessor, GPT-3.5. According to the report, GPT-4 exhibits human-level performance on various professional and academic benchmarks, and specifically, in a translated exam benchmark, it outperforms previous models in 24 out of 26 languages tested [26]. This indicated that GPT-4 has a strong grasp of not only English but a wide range of languages, including some low-resource ones, likely due to the diversification of its training data and increased model capacity. The technical report also discussed efforts to make GPT-4 outputs more reliable and aligned via fine-tuning and reinforcement learning. Significance: GPT-4 represented the state of the art in LLMs, and its confirmed prowess in multilingual understanding meant that general AI systems were becoming direct competitors to specialized MT systems. The report reinforced an emerging view in 2023: organizations planning for the future of localization must consider LLMs as part of the toolkit - whether for generating draft translations, performing guality assessment, or enabling multilingual chatbots. However, the GPT-4 report also cautioned about issues like hallucinations and inconsistencies, reminding stakeholders that expert human oversight remains essential when using these models for critical translation tasks.

Conclusion

The last decade (2014–2024) has seen unprecedented progress at the intersection of AI and GILT. We began with the introduction of neural networks into MT, which rapidly evolved from basic RNN-based seq2seq models to complex attention-based Transformers, culminating in massive multilingual systems that aim to serve all languages. Translation quality on major benchmarks has improved to the point of nearing human performance in certain settings, driven by innovations like attention mechanisms [6], subword modeling [8], and effective use of monolingual data [9]. At the same time, the rise of large language models (GPT-3 [20], GPT-4 [26]) has expanded the scope of what AI can do – translating not as a narrow task-specific system, but as one of many capabilities of a general language assistant [2] [25].

Throughout these developments, the role of the human linguist and the importance of human factors have become ever more apparent. Studies on post-editing and interactive MT remind us that technology must serve its users: interface design, adaptivity, and translator training are crucial to achieving real productivity gains [5] [24]. The notion of "human parity" in translation, while tantalizing, has been nuanced by research that digs deeper into evaluation conditions [15]. Rather than viewing AI as replacing human translators, the field has increasingly moved toward viewing it as a powerful tool that – when properly aligned and integrated – amplifies human capabilities.

As of 2024, AI in GILT stands at a crossroads. Large multilingual models and LLMs offer the possibility of translating and localizing content in hundreds of languages instantly, potentially breaking remaining language barriers. Yet challenges remain: ensuring factual and cultural accuracy, handling low-resource languages with limited data, maintaining translation quality in specialized domains, and doing all this in a way that is ethical, fair, and preserves linguistic diversity. The academic works reviewed in this monograph collectively illustrate a trajectory of remarkable achievements addressing these challenges one by one. They also set the stage for the next era of research – one that may focus on fine-grained control of translation style and tone, deeper contextual and multimodal understanding (e.g., translating not just text but integrating images or video context), and robust evaluation metrics that keep AI developers honest about claiming "human quality."

In summary, the period 2014–2024 transformed AI in the GILT industry from a niche endeavor into a central pillar of how global content is created and consumed. Ongoing collaboration between computational researchers, linguists, and localization professionals – as exemplified by the papers discussed – will be critical to ensure that the next decade's innovations continue to advance both the science of translation and the art of global communication.

References

1. Cadieux, P., & Esselink, B. (2002). GILT: Globalization, Internationalization, Localization, Translation. Retrieved from <u>https://www.i18n.ca/workshops/pdf/GILT.pdf</u>

2. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gómez, A. N., Kaiser, Ł., & Polosukhin, I. (2017). Attention is all you need. In Advances in Neural Information Processing Systems 30 (pp. 5998–6008).

3. Cho, K., van Merriënboer, B., Gulcehre, C., Bahdanau, D., Bougares, F., Schwenk, H., & Bengio, Y. (2014). Learning phrase representations using RNN Encoder–Decoder for statistical machine translation. In Proceedings of EMNLP 2014 (pp. 1724–1734).

4. Sutskever, I., Vinyals, O., & Le, Q. V. (2014). Sequence to sequence learning with neural networks. In Advances in Neural Information Processing Systems 27 (pp. 3104–3112).

5. Green, S., Wang, S., Chuang, J., Heer, J., Schuster, S., & Manning, C. D. (2014). Human effort and machine learnability in computer-aided translation. In Proceedings of EMNLP 2014 (pp. 1289–1299).

6. Bahdanau, D., Cho, K., & Bengio, Y. (2015). Neural machine translation by jointly learning to align and translate. In 3rd International Conference on Learning Representations (ICLR).

7. Zaretskaya, A., Corpas Pastor, G., & Seghiri, M. (2015). Integration of machine translation in CAT tools: State of the art, evaluation and user attitudes. SKASE Journal of Translation and Interpretation, 8(1), 76–88.

8. Sennrich, R., Haddow, B., & Birch, A. (2016). Neural machine translation of rare words with subword units. In Proceedings of ACL 2016 (pp. 1715–1725).

9. Sennrich, R., Haddow, B., & Birch, A. (2016). Improving neural machine translation models with monolingual data. In Proceedings of ACL 2016 (pp. 86–96).

10. Wu, Y., Schuster, M., Chen, Z., Le, Q. V., Macherey, W., et al. (2016). Google's neural machine translation system: Bridging the gap between human and machine translation. arXiv:1609.08144.

References (continued)

11. Johnson, M., Schuster, M., Le, Q. V., Krikun, M., Wu, Y., Chen, Z., ... & Dean, J. (2017). Google's multilingual neural machine translation system: Enabling zero-shot translation. Transactions of the ACL, 5, 339–351.

12. Koehn, P., & Knowles, R. (2017). Six challenges for neural machine translation. In Proceedings of the First Workshop on Neural Machine Translation (pp. 28–39).

13. Moorkens, J., & O'Brien, S. (2017). Assessing user interface needs of post-editors of machine translation. In F. Sharmin & D. Kenny (Eds.), Human Issues in Translation Technology (pp. 127–148). Routledge.

14. Hassan, H., Aue, A., Chen, C., Chowdhary, V., Clark, J., ... & Liu, T.-Y. (2018). Achieving human parity on automatic Chinese to English news translation. arXiv:1803.05567.

15. Toral, A., Castilho, S., Hu, K., & Way, A. (2018). Attaining the unattainable? Reassessing claims of human parity in neural machine translation. In Proceedings of the 3rd Conference on Machine Translation (WMT 2018) (pp. 113–123).

16. Lample, G., Conneau, A., Denoyer, L., & Ranzato, M. (2018). Unsupervised machine translation using monolingual corpora only. In Proceedings of ICLR 2018.

17. Devlin, J., Chang, M.-W., Lee, K., & Toutanova, K. (2019). BERT: Pre-training of deep bidirectional transformers for language understanding. In Proceedings of NAACL-HLT 2019 (pp. 4171–4186).

18. Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., & Sutskever, I. (2019). Language models are unsupervised multitask learners. OpenAI Technical Report.

19. Daems, J., & Macken, L. (2019). Interactive adaptive SMT versus interactive adaptive NMT: A user experience evaluation. Machine Translation, 33(1-2), 117–134.

20. Brown, T. B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., ... & Amodei, D. (2020). Language models are few-shot learners. In Advances in NeurIPS 33 (pp. 1877–1901).

References (continued)

21. Fan, A., Bhosale, S., Schwenk, H., Ma, Z., El-Kishky, A., ... & Joulin, A. (2021). Beyond English-centric multilingual machine translation. Journal of Machine Learning Research, 22(107), 1–48.

22. Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C., ... & Christiano, P. (2022). Training language models to follow instructions with human feedback. arXiv:2203.02155.

23. NLLB Team, Costa-jussà, M. R., Cross, J., Çelebi, O., Elbayad, M., ... & Wang, J. (2022). No language left behind: Scaling human-centered machine translation. arXiv:2207.04672.

24. Martikainen, H. (2022). Ghosts in the machine: Can adaptive MT help reclaim a place for the human in the loop? [Preprint]. hal-03548696.

25. Yan, J., Yan, P., Chen, Y., Li, J., Zhu, X., & Zhang, Y. (2023). GPT-4 vs. human translators: A comprehensive evaluation of translation quality across languages, domains, and expertise levels. arXiv:2407.03658.

26. OpenAI. (2023). GPT-4 Technical Report. arXiv:2303.08774.

GILT Research Centre (Globalisation, Internationalisation, Localisation and Translation)

2024