## "Critical Summary of Liu et al., 2020: Multilingual Denoising Pre-training for Neural Machine Translation"

This paper presents **mBART**, a sequence-to-sequence multilingual pre-trained model designed to enhance Neural Machine Translation (NMT) tasks, particularly for low-resource and unsupervised scenarios. In contrast to earlier methods like XLM-R or mBERT, which focus solely on encoder pre-training or employ sentence-level Masked Language Modeling (MLM) objectives, mBART pre-trains both the encoder and decoder simultaneously using a denoising autoencoder approach. The authors illustrate mBART's applicability across sentence- and document-level translation tasks and highlight its capability to generalize to unseen languages, achieving notable improvements in diverse contexts. The mBART architecture is based on a Transformer sequence-to-sequence model with 12 layers each in the encoder and decoder, amounting to approximately 680M parameters. Pre-training utilizes CC25, a rebalanced dataset comprising monolingual data from 25 languages, ensuring representation across a wide range of linguistic characteristics. The pre-training process employs a denoising objective where input text is corrupted through *span masking* (randomly masking 35% of tokens) and *sentence permutation* (shuffling sentence order) before reconstruction. A SentencePiece tokenizer, trained on data from 100 languages, ensures compatibility with both pre-training and fine-tuning tasks and supports generalization to languages not present during pre-training. The model and procedure were evaluated for, Sentence-level Translation: Significant gains were observed for low-resource languages, such as a 12 BLEU-point improvement in English-to-Vietnamese translation. Document-level Translation: The model, trained on WMT19 data and evaluated on TED datasets, demonstrated enhanced coherence and fluency over non-pre-trained baselines. In Unsupervised NMT mBART was tested on tasks where no parallel bi-text data existed for target language pairs, achieving strong performance in back-translation and transfer tasks. Promising results were also noted when combining back-translation with language transfer.

**Strengths and Limitations:** mBART's strengths include its robust performance in low/medium-resource languages, with up to 12 BLEU-point improvements, and its ability to enable transfer learning for unseen languages, enhancing its versatility. In unsupervised settings, mBART showed consistent improvements, achieving the first non-degenerate results for less related language pairs, such as a 9.5 BLEU-point gain in Nepali-to-English translation. Additionally, the joint pre-training of the encoder and decoder aligns naturally with sequence-to-sequence tasks, unlike earlier models optimized for classification tasks. However, the model's limitations include its dependence on extensive monolingual corpora and significant computational resources, requiring 256 GPUs over 2.5 weeks for training. Reproducibility poses another challenge, as re-training can yield slightly variable fine-tuning performance. Furthermore, for high-resource language pairs, mBART shows minimal improvement, with pre-training even slightly degrading performance when datasets exceed 10M sentence pairs. The study also leaves unexplored the potential of alternative noise functions, which may further optimize pre-training.

## Differences between mBART and mBERT/XLM-R:

- Architecture: mBERT/XLM-R use encoder-only architectures, while mBART includes both an encoder and a decoder, optimized for sequence-to-sequence tasks.
- Pre-training Task/Noise Function: mBERT and XLM-R rely on masked language modeling, whereas mBART uses both span masking and sentence permutation (shuffling sentence order) to reconstruct full sequences.
- Pre-training Data: mBERT/XLM-R train on multilingual Wikipedia (or CC100 for XLM-R), while mBART pre-trains on rebalanced CC25 (25 languages) data to ensure better representation for low-resource language.

**Tokenizer Design:** The SentencePiece tokenizer was trained on 100 languages, extending beyond the 25 languages used in mBART25. This decision facilitates generalization to unseen languages during fine-tuning. This approach ensures that token embeddings are shared across related languages, enabling zero-shot, transfer learning (broader cross-lingual transfer) but may introduce inefficiency with rarely used tokens (large vocabulary).

**High- vs. Low-Resource Language Benefits:** High-resource languages show limited improvements, with pre-training even slightly degrading performance when parallel data exceeds 25M sentences because their substantial parallel data allows fine-tuning to wash out/override pre-trained weights. In contrast, low-resource languages benefit significantly from mBART pre-training, as it provides crucial linguistic priors lacking in their limited training data.

**Conclusion:** mBART is a significant advancement in multilingual pre-training, demonstrating improvements across diverse translation settings. While computational costs and high-resource performance pose challenges, its impact on low-resource and unseen language tasks highlights its potential for expanding multilingual NLP capabilities. Future research could explore efficiency improvements and scaling to additional languages.

**Unclear points:** The details about denoising functionality and its architecture could be helpful to better understand. Additionally, insights into how the model balances the multilingual data during training to prevent dominance by high-resource languages would be beneficial.