

1 Appendix-A: Model and Dataset

2 In this section, we provide details on the datasets, model architectures and settings used for full
3 precision fine-tuning to replicate the baseline accuracy reported in publications.

4 1.1 BERT-base/SQuAD1.1

5 **SQuAD1.1** is the Stanford Question Answering dataset [1]. It is a reading comprehension dataset,
6 consisting of a collection of 100k question and answer pairs from reading passages, where the answer
7 to every question is a segment of corresponding passage. The task is to predict the answer text span
8 in a passage.

9 **BERT-base** model [2] is a transformer model pre-trained on large corpus of English data, i.e.
10 BooksCorpus [3] with 800M words and English Wikipedia of 2500M words. The model was pre-
11 trained in a self-supervised fashion using a masked language modeling (MLM) procedure. The
12 BERT-base model consists of an input embedding, 12 transformer blocks and an output linear layer,
13 with the total parameters of 110M. The input embedding is a sum of token embeddings, segmentation
14 embeddings and the position embeddings. Each transformer block contains 12 self-attention heads
15 and a hidden size of 768.

16 For fine-tuning BERT-base on SQuAD1.1 downstream task, we use batch size of 12 and sequence
17 length of 384. The experiments are performed on 4 V100 GPUs with per-gpu batch size of 3. For the
18 full precision fine-tuning baseline, we follow the fine-tuning strategy from [2] and use the AdamW
19 optimizer with a learning rate of $3e-5$ with linear decay. The model is fine-tuned for 2 epochs with
20 a dropout probability of 0.1. We obtain a baseline F1 score of 88.69 which closely matches the F1
21 score (88.50) published in [2]. Fig. 1a shows the convergence curve of the full precision fine-tuning
22 baseline.

23 1.2 Wav2vec2.0/Librispeech

24 **Librispeech** is a corpus of English speech for automatic speech recognition (ASR) task [4]. It
25 contains 1000 hours of speech sampled at 16 kHz. The training data is split into 3 partitions of 100hr,
26 360hr and 500hr sets with 'clean' and 'other' categories. In this work, we use the 100hr-clean data
27 for downstream task and the clean validation subset for evaluation.

28 **Wav2vec2.0-large** is speech model pre-trained on the audio data from LibriVox (LV-60k) [5] in a
29 self-supervised manner [6]. In this work, we use the Wav2vec2.0-large model which has a large
30 transformer backbone with 24 transformer blocks. The hidden dimension, inner dimension and
31 number of attention heads in each transformer block are 1024, 4096 and 16, respectively. Before fed
32 into transformer backbone, the raw waveform is encoded through multiple 1d-convolution layers
33 followed by layer normalization and GeLU activations.

34 The pre-trained model is fine-tuned on Librispeech's 100 hour clean subset using standard Con-
35 nectionist Temporal Classification (CTC) loss. We follow the implementation and settings from
36 HuggingFace Transformer [7] for the fine-tuning. Specifically, for full precision fine-tuning baseline,
37 we use the AdamW optimizer with betas=(0.9,0.999) and a learning rate of $3e-4$. The learning rate
38 decays linearly after 500 warm-up steps. We use 8 V100 GPUs to tune the model for 3 epochs with a
39 total batch size of 32. We achieve baseline Word Error Rate (WER) of 4.20 %, matching the result
40 provided by HuggingFace Transformer (4.2 %). Fig. 1b shows the convergence curve of full precision
41 fine-tuning baseline.

42 1.3 ViT/ImageNet1k

43 **ImageNet1k** [8] is an image classification benchmark which consists of 1000-categories of objects
44 with over 1.2M training and 50K validation images.

45 **ViT-base** model is a BERT-like transformer encoder model taking images as the input for image
46 classification tasks. The input images are split into fixed-sized patches of 16x16 and linearly
47 embedded. The ViT-base model has 12 transformer blocks with 12 attention heads and a hidden
48 dimension of 768 [9]. In this paper, we use the pre-trained model that is only pre-trained on
49 ImageNet21k [8] and then fine-tune it on ImageNet1k for downstream image classification. We

Table 1: Coefficients for SAWB+.

Precision	C1	C2
INT4	12.68,	-12.80
INT8	31.76,	-35.04

50 use a resolution of 384x384 for fine-tuning, following the original settings of [9]. The optimizer
 51 is SGD with a learning rate of 0.01. We tune the model for 8 epochs using a Cosine learning rate
 52 schedule, gradient clipping of 1.0, and batch size of 512 on 32 V100 GPUs. With these settings, we
 53 achieve accuracy of 84.12 which matches the accuracy (83.97) published in [9]. Fig. 1c shows the
 54 convergence curve of full precision fine-tuning baseline.

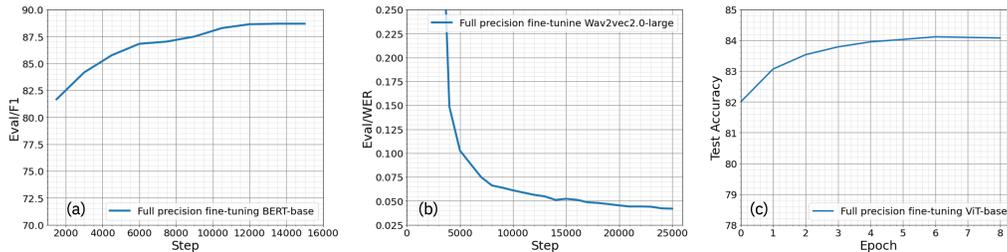


Figure 1: Convergence curves of full precision fine-tuning of (a) BERT-base on SQuAD1.1; (b) Wav2vec2.0-large on Librispeech; and (c) ViT-base on ImageNet1k.

55 2 Appendix-B: Quantization- and Sparsity-aware Fine-tuning Settings

56 In this section, we provide details on the implementation of quantization and pruning operations, as
 57 well as the hyper-parameters used for the fine-tuning of deep compressed models.

58 2.1 Quantization/pruning implementation

59 We implement the quantization and pruning in PyTorch framework. For each linear module or
 60 batch-matrix-matrix multiplication (bmm) operation, we insert quantization operations to quantize
 61 both the activation and weight. Fig. 2 shows screenshot examples of quantized modules with inserted
 62 quantization operations from the graph of quantized BERT-base model. For linear modules, both
 63 input activation and weight are quantized, as shown in Fig. 2(a) a query layer in self-attention and
 64 (c) an intermediate-dense layer in the feed-forward network (FFN); while, for bmm operations, both
 65 input activations are quantized as shown in Fig. 2(b).

66 Fig. 3 presents a toy example showing a deep compressed linear module, i.e. QLinear, running a
 67 forward pass with a random input. The linear layer is quantized in 4-bit for both weight and activation
 68 using SAWB+ and PACT quantizers, respectively. The weight is further pruned with 50% sparsity
 69 using the a fine-grain group of 4 as discussed in section 2.2. The printout shows the pruning mask
 70 tensor, pruned weight tensor and quantized weight tensor computed during the forward pass.

71 2.2 Fine-tune setting

72 We use a common setting for all three models and benchmarks. Full precision fine-tuned models are
 73 used for the initialization of INT8, sparse INT8 and INT4 models. For the sparse INT4 model, we
 74 use a sparse INT8 model for initialization as explained in section 2.3. Fig. 4 shows a schematic of
 75 the fine-tuning procedures. The SAWB+ quantizer is used for weight quantization for all models as
 76 discussed in section 2.1.1. The coefficients used in SAWB+ are listed in Table 1. The MinMax or
 77 PACT quantizer is used for the activation quantization. For PACT quantizer (discussed in section
 78 2.1.2), three hyper-parameters are used to train α and α_n parameters, i.e. initiation in percentile,

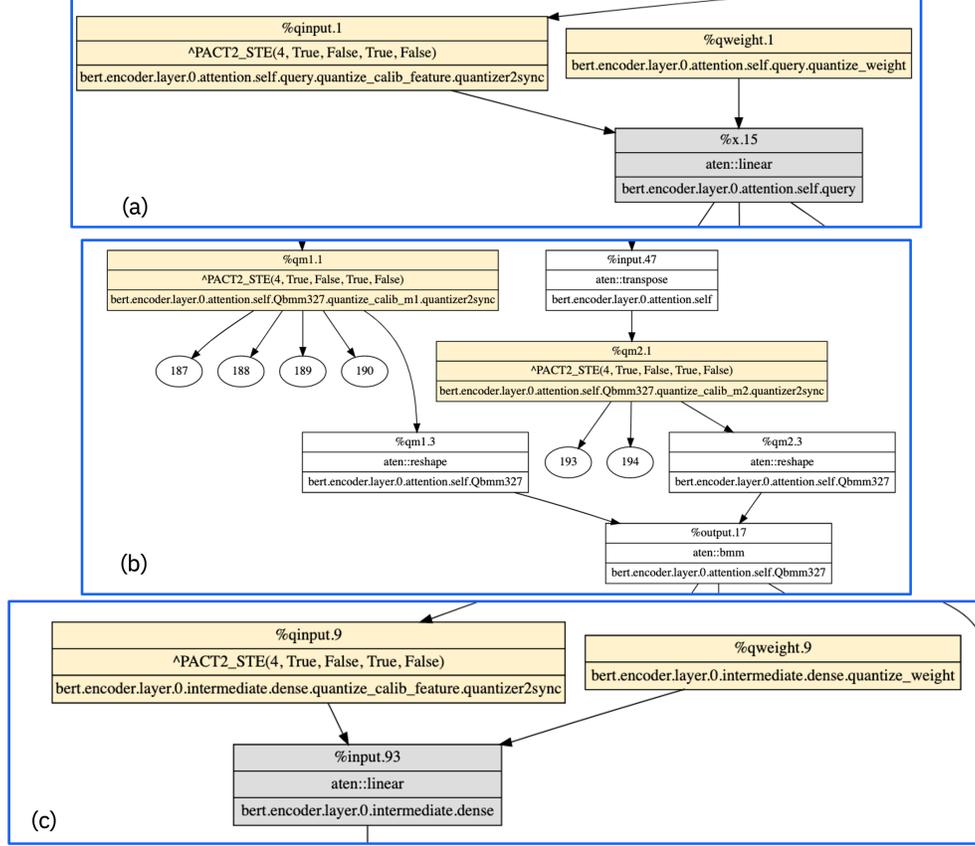


Figure 2: a) Screenshot examples of a quantized graph with implemented weight and activation quantization operations, for (a) a query linear layer; (b) a bmm operation for attention computation; and (c) an intermediate dense linear layer in FFN from layer0 (the first transformer block) of the deep compressed BERT-base model.

Table 2: Quantization/Sparsity-aware fine-tuning setting for BERT-base on SQuAD1.1. Sp is short for sparsity.

Precision Sparsity	Weight Quantizer	Activation Quantizer	Initialization Model	Percentile (%)	α_{lr}	α_{decay}	Dropout
INT8	SAWB+	MinMax	FP32	–	–	–	0.2
INT8+50%Sp	SAWB+	MinMax	FP32	–	–	–	0.2
INT4	SAWB+	PACT	FP32	99	1e-3	1e-3	0.2.
INT4+50%Sp	SAWB+	PACT	INT8+50%Sp	99	1e-3	1e-3	Scheduled

79 learning rate (α_{lr}) and L2 decay (α_{decay}). The detailed settings used for three benchmarks are as
80 follows.

81 Table 2 lists the settings for BERT-base/SQuAD1.1 benchmark. We use the same baseline optimization
82 methods as described in Appendix 1.1, except that the compressed models are fine-tuned for 4 epochs
83 with a larger dropout (0.2) or a scheduled dropout as introduced in section 2.3.3. Fig. 5 shows the
84 convergence curves of the deep compressed models.

85 Table 3 lists the settings for Wav2vec2.0-large/Librispeech benchmark. We use the same baseline
86 optimization methods as described in Appendix 1.2, except that the compressed models are fine-tuned
87 for 6 epochs. Fig. 6 shows the convergence curves of the deep compressed models.

88 Table 4 lists the settings for ViT-base/ImageNet1k benchmark. For INT8/4 models without pruning,
89 we use the same baseline optimization methods as described in Appendix 1.3. For sparse INT8/4

```

1 #a toy example to quantize and prune one linear layer with dummy inputs
2 model_q = QLinear(in_features=64, out_features=4, num_bits_feature=4, num_bits_weight=4, \
3                 p_group=4, p_ratio = 0.5, qa_mode='pact', qw_mode='sawb+')
4 #get pruning mask
5 model_q.get_mask()
6 output = model_q(torch.rand(1, 4, 64))

Weights[:, :8]:
tensor([[ 1.1034e-01,  7.0334e-02, -9.7582e-02,  4.9689e-02, -8.1956e-07,
        -9.7793e-02, -8.3745e-06, -1.0821e-01],
        [-1.1906e-01,  8.3641e-02,  3.2591e-02,  3.1367e-02, -1.5860e-02,
         6.3644e-02, -1.1322e-01, -1.1617e-01],
        [ 2.6805e-02, -3.2602e-02,  1.0721e-01, -3.6630e-02, -1.1574e-01,
         7.5815e-02, -5.7634e-02, -1.1444e-02],
        [ 1.1820e-01,  5.0626e-02, -7.3768e-02,  6.0336e-02, -8.2874e-02,
        -6.0145e-02,  1.0179e-01,  9.6717e-02]])

Mask[:, :8]:
tensor([[1., 0., 1., 0., 0., 1., 0., 1.],
        [1., 1., 0., 0., 0., 0., 1., 1.],
        [0., 0., 1., 1., 1., 1., 0., 0.],
        [1., 0., 1., 0., 0., 0., 1., 1.]])

Pruned weight[:, :8]:
tensor([[ 0.1103,  0.0000, -0.0976,  0.0000, -0.0000, -0.0978, -0.0000, -0.1082],
        [-0.1191,  0.0836,  0.0000,  0.0000, -0.0000,  0.0000, -0.1132, -0.1162],
        [ 0.0000, -0.0000,  0.1072, -0.0366, -0.1157,  0.0758, -0.0000, -0.0000],
        [ 0.1182,  0.0000, -0.0738,  0.0000, -0.0000, -0.0000,  0.1018,  0.0967]])

Quantized weight[:, :8]:
tensor([[ 0.1067,  0.0000, -0.1067,  0.0000,  0.0000, -0.1067,  0.0000, -0.1067],
        [-0.1067,  0.0711,  0.0000,  0.0000,  0.0000,  0.0000, -0.1067, -0.1067],
        [ 0.0000,  0.0000,  0.1067, -0.0356, -0.1067,  0.0711,  0.0000,  0.0000],
        [ 0.1067,  0.0000, -0.0711,  0.0000,  0.0000,  0.0000,  0.1067,  0.1067]])

```

Figure 3: a) A toy example of a quantized and pruned linear module running a forward pass with a random input. The printout shows the pruning mask, pruned weight and quantized weight tensors.

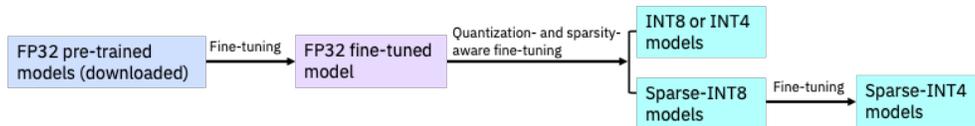


Figure 4: A schematic of the fine-tuning procedures leading to deep compressed models. The quantization- and sparsity- aware fine-tuning are initialized by FP32 fine-tuned models to obtain the INT8, INT4 or sparse-INT8 models. The sparse-INT8 models are further fine-tuned to get the sparse-INT4 models.

90 models, we tune the model for 16 epochs with starting learning rate of 0.05, keeping the rest of hyper-
91 parameters the same as the baseline. Fig. 7 shows the convergence curves of the deep compressed
92 models.

93 3 Appendix-C: Broader Impact

94 Dedicated hardware accelerators for DNN inference, including CPUs, GPUs, TPUs and other AI
95 platforms, have powered the deployment of machine learning for real-life applications in both cloud
96 and edge devices. Reduced precision innovations (FP16, FP8 and INT8), together with sparsity,
97 have recently improved the capability of these accelerators by 4-8 \times and have dramatically improved

Table 3: Quantization/Sparsity-aware fine-tuning setting for Wav2vec2.0-large on Librispeech. Sp is short for sparsity.

Precision Sparsity	Weight Quantizer	Activation Quantizer	Initialization Model	Percentile (%)	α_{lr}	α_{decay}
INT8	SAWB+	MinMax	FP32	–	–	–
INT8+50%Sp	SAWB+	MinMax	FP32	–	–	–
INT4	SAWB+	PACT	FP32	max	1e-2	7e-3
INT4+50%Sp	SAWB+	PACT	INT8+50%Sp	99.9	1e-2	3e-2

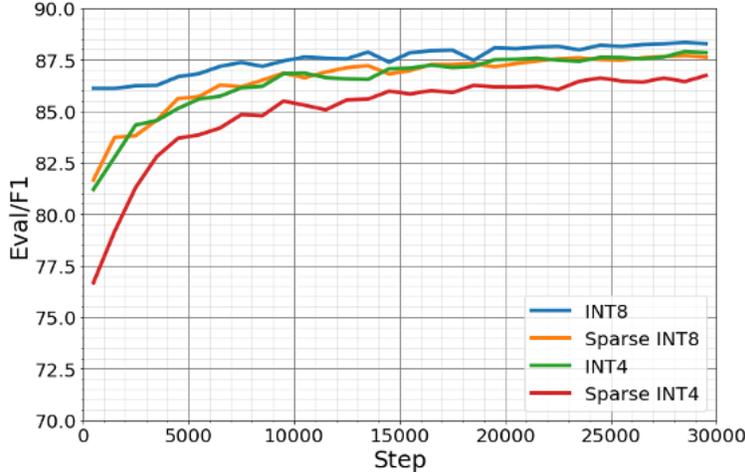


Figure 5: Convergence curves of INT8, sparse INT8, INT4 and sparse INT4 BERT-base models on SQuAD1.1.

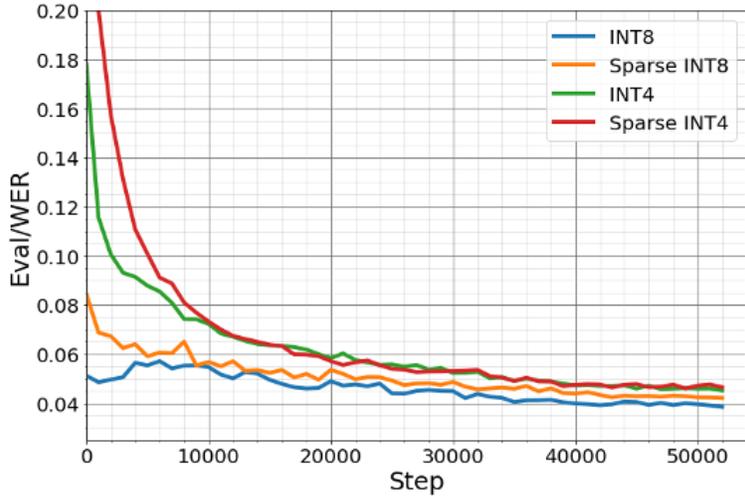


Figure 6: Convergence curves of INT8, sparse INT8, INT4 and sparse INT4 Wav2vec2.0-large models on Librispeech.

98 energy cost and carbon emissions. Although pre-trained transformers have unlocked the power of
 99 transfer learning and are leading to breakthroughs in multiple application domains, the architecture
 100 is too complex for many production systems, such as those for edge-computing inference. There
 101 are many ongoing efforts to reduce the size of these models while retaining model performance and
 102 transferability. Deep compression of transformers, which is presented in this work, aims to push
 103 this front aggressively to enable faster and cheaper inference systems for a wide spectrum of deep
 104 learning models and domains. We believe that sparse 4-bit inference solutions can accelerate ML
 105 deployment and provide significant cost and energy savings for corporations and research institutes
 106 — in addition to helping reduce the carbon / climate impact of AI inference. By improving power
 107 efficiency by about $4\times$ over current transformers running in FP16 (and $8\times$ vs. default FP32 designs),
 108 the carbon footprint for predicting with large DNN models can be significantly reduced [10].

109 The reduction in computational energy and memory footprint could also enable the inference of
 110 large transformer models to be carried out on edge devices (mobile platforms, health care devices,
 111 security cameras, consumer drones, etc.). This, in turn, could alleviate security and privacy concerns
 112 of sending data back to the Cloud for prediction tasks.

Table 4: Quantization/Sparsity-aware fine-tuning setting for ViT-base on ImageNet1k. Sp is short for sparsity.

Precision Sparsity	Weight Quantizer	Activation Quantizer	Initialization Model	Percentile (%)	α_{lr}	α_{decay}	Epoch
INT8	SAWB+	MinMax	FP32	–	–	–	8
INT8+50%Sp	SAWB+	MinMax	FP32	–	–	–	8
INT4	SAWB+	PACT	FP32	99.9	1e-2	1e-5	16
INT4+50%Sp	SAWB+	PACT	INT8+50%Sp	99.9	1e-2	1e-6	16

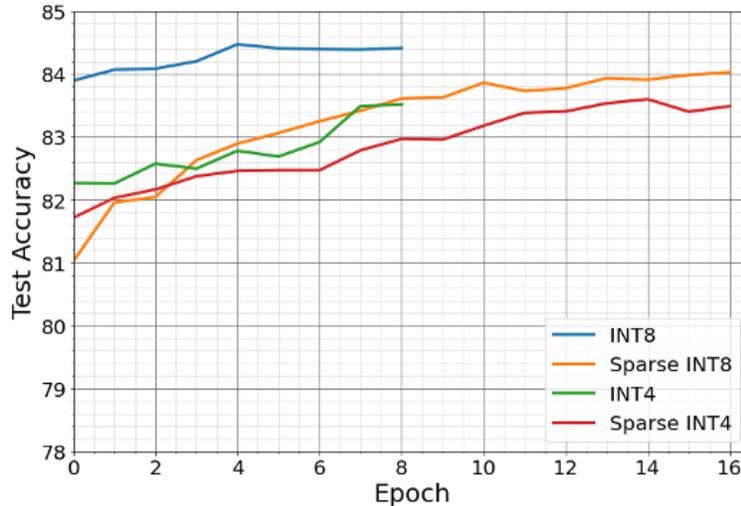


Figure 7: Convergence curves of INT8, sparse INT8, INT4 and sparse INT4 ViT-base models on ImageNet1k.

113 We would also like to emphasize that, although we have shown promising results and limited accuracy
 114 loss in comparison to FP32 downstream tasks, deep compressed transformer models using our
 115 solutions could still be subject to unexpected instabilities. This may necessitate a careful examination
 116 of these optimization techniques and numerical formats over a wider range of models and perfected
 117 alongside the development of ML model research. The risk of using deeply compressed transformer
 118 models in real inference applications is most likely higher than full precision dense models and thus
 119 requires task-specific robustness studies to prepare these models against adversarial attacks. More
 120 work is also needed to assess the impact of deeply compressed models in fairness and explainability.

121 References

- 122 [1] Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. Squad: 100, 000+
 123 questions for machine comprehension of text. *CoRR*, abs/1606.05250, 2016. URL <http://arxiv.org/abs/1606.05250>.
 124
- 125 [2] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of
 126 deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*,
 127 2018.
- 128 [3] Yukun Zhu, Ryan Kiros, Richard Zemel, Ruslan Salakhutdinov, Raquel Urtasun, Antonio
 129 Torralba, and Sanja Fidler. Aligning books and movies: Towards story-like visual explanations
 130 by watching movies and reading books, 2015.
- 131 [4] Vassil Panayotov, Guoguo Chen, Daniel Povey, and Sanjeev Khudanpur. Librispeech: An
 132 asr corpus based on public domain audio books. In *2015 IEEE International Conference on*
 133 *Acoustics, Speech and Signal Processing (ICASSP)*, pages 5206–5210, 2015. doi: 10.1109/
 134 ICASSP.2015.7178964.

- 135 [5] J. Kahn, M. Riviere, W. Zheng, E. Kharitonov, Q. Xu, P.E. Mazare, J. Karadayi, V. Liptchin-
136 sky, R. Collobert, C. Fuegen, T. Likhomanenko, G. Synnaeve, A. Joulin, A. Mohamed,
137 and E. Dupoux. Libri-light: A benchmark for ASR with limited or no supervision. In
138 *ICASSP 2020 - 2020 IEEE International Conference on Acoustics, Speech and Signal*
139 *Processing (ICASSP)*. IEEE, may 2020. doi: 10.1109/icassp40776.2020.9052942. URL
140 <https://doi.org/10.1109%2Ficassp40776.2020.9052942>.
- 141 [6] Alexei Baevski, Henry Zhou, Abdelrahman Mohamed, and Michael Auli. wav2vec 2.0: A frame-
142 work for self-supervised learning of speech representations. *arXiv preprint arXiv:2006.11477*,
143 2020.
- 144 [7] et al. Wolf, Thomas. Huggingface’s transformers: State-of-the-art natural language processing.
145 *arXiv preprint arXiv:1910.03771*, 2019.
- 146 [8] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-
147 scale hierarchical image database. In *2009 IEEE Conference on Computer Vision and Pattern*
148 *Recognition*, pages 248–255, 2009. doi: 10.1109/CVPR.2009.5206848.
- 149 [9] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai,
150 Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al.
151 An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint*
152 *arXiv:2010.11929*, 2020.
- 153 [10] Emma Strubell, Ananya Ganesh, and Andrew McCallum. Energy and policy considerations for
154 deep learning in nlp. *arXiv preprint arXiv:1906.02243*, 2019.