# The Privileged Students: On the Value of Initialization in Multilingual Knowledge Distillation

Anonymous ACL submission

#### Abstract

Knowledge distillation (KD) has proven to be a successful strategy to improve the performance of smaller models in many NLP tasks. However, most of the work in KD only explores monolingual scenarios. In this paper, we investigate the value of KD in multilingual settings. We find the significance of KD and model initialization by analyzing how well the student model acquires multilingual knowledge from the teacher model. Our findings show that model initialization using copy-weight from the fine-tuned teacher contributes the most compared to the distillation process itself across various multilingual settings. Furthermore, we demonstrate that efficient weight initialization preserves multilingual capabilities even in lowresource scenarios, enabling the student models to perform to languages unseen during distillation.

### 1 Introduction

001

006

800

011

012

014

017

018

023

037

Recent work has shown that knowledge distillation (KD) is effective in distilling multilingual language models (Hui et al., 2024; Grattafiori et al., 2024; Ansell et al., 2023; Team et al., 2024). However, how multilingual capabilities are preserved during the distillation process remains understudied.

Understanding the mechanisms of multilingual KD is crucial for improving best practices. Some of the main challenges in multilingual NLP are the scarcity of data and limited computational resources (Conneau et al., 2020). By uncovering how cross-lingual information is transferred during multilingual KD, we can train lightweight multilingual models more effectively, reducing the need for extensive data, thus enhancing accessibility.

The typical KD process involves two stages pattern (Jiao et al., 2020): initialization, where the student model receives knowledge from the teacher model, and fine-tuning, where the student learns downstream tasks through distillation loss. In this work, we analyze which component contributes more to preserving multilingual capabilities in resource-constrained settings. 041

042

043

044

045

046

047

051

056

057

060

061

062

063

064

065

066

067

068

069

071

072

073

074

075

076

077

We are also intrigued whether the teacher's cross-lingual generalization abilities transfer to the student model through this process. Understanding this transfer capability is important as it could enable more efficient deployment of multilingual models in resource-constrained environments. Moreover, we examine how efficiently knowledge is transferred across different training data sizes. Achieving multilingual capabilities through simple weight copying, rather than full pre-training, would represent a significant advancement in efficiency.

Our findings indicate that weight copying plays a more important role, even more than the KD loss itself. Specifically, we observe that models with copy-weight initialization outperform freshly initialized models trained using a KD loss. Moreover, properly initialized student models demonstrate the ability to generalize to unseen languages during the distillation process—an outcome unattainable when relying solely on the KD loss (Ansell et al., 2023).

Based on the above motivations, our contributions to this work are as follows.

- 1. We found that model initialization through weight copying plays a more crucial role than the distillation process itself in preserving multilingual performance.
- 2. We show that distilled model can generalize towards languages unseen during distillation training, as long as the model is initialized properly.
- 3. We identify that the weight-copy approach leads to more efficient training, exhibiting faster convergence across varying amounts of training data.



Figure 1: We investigate the multilingual capability of distilled model with copy-weight as the initialization. We found that, it moderately preserves multilingual capability of the teacher, which happily make the learning data efficient and faster convergence.

#### 2 Background: Knowledge Distillation

079

084

091

096

099

100

101

102

103

104

105

108

109

110

111

Knowledge distillation (KD) is a technique used to transfer knowledge from a large, trained teacher model to a smaller student model. This process aims to retain the large model's performance while reducing the computational cost during inference. KD involves training the student model to mimic the outputs of the teacher model, using a combination of the teacher's soft target outputs and the ground truth labels as the learning objective.

Jiao et al. (2020) introduces a two-step distillation process. First, the student model is pre-trained on a large corpus to acquire good initialization suitable for the next step. Afterward, the model is fine-tuned for the desired tasks. Unfortunately, The former approach requires substantial data and computational resources, making it challenging to implement with limited resources.

In this work, we explore the impact of these components of knowledge distillation in multilingual settings. Instead of the extensive pre-training step used in Jiao et al., 2020; Ansell et al., 2023, we investigated a simpler and more efficient initialization approach by copying the weights from the teacher model to the student model proposed by DistilBERT (Sanh et al., 2020). By copying some of the layers, They may preserve multilingual information that is beneficial for the smaller model (See Figure 1 for the illustration). We elucidate the possibility that this approach has already transferred the multilingual transfer even without further pre-training.

The underlying mechanism of KD and the copy

weight are elaborated in this section.

#### 2.1 Distillation Architecture

We utilize KD, comprised of a teacher T and a student S model. The student model always has fewer layers than the teacher model. We follow TinyBERT (Jiao et al., 2020)'s objective loss and architecture. The loss of the KD comprises embedding loss  $\mathcal{L}_{embd}$ , hidden-layer loss  $\mathcal{L}_{hidn}$ , attention loss  $\mathcal{L}_{att}$ , and prediction-layer loss  $\mathcal{L}_{pred}$ . These objective functions can be formulated as follows:

$$\mathcal{L}_{att} = \frac{1}{l} \frac{1}{h} \sum_{i=1}^{l} \sum_{j=1}^{h} \text{MSE}(A_{S}^{i}, A_{T}^{k})$$
(1)

112

113

114

115

116

117

118

119

121

122

125

126

127

128

129

130

131

132

133

134

$$\mathcal{L}_{hid} = \frac{1}{l} \sum_{i=1}^{l} \text{MSE}(W \cdot H_S^i, H_T^k)$$
(2)

$$\mathcal{L}_{embd} = \mathrm{MSE}(W \cdot E_S, E_T) \tag{3}$$

$$\mathcal{L}_{pred} = \text{MSE}(z_S, z_T) \tag{4}$$

Where A, H, E, and z are the values of the attention outputs, hidden layers' outputs, embedding layer's outputs, and the logits, respectively, for the teacher T or student S models. Indices of model layers and attention heads are denoted as l and h. If the student model's hidden unit dimension is smaller than the teacher's, we leverage a projection weight W to match the hidden unit dimension. Otherwise, W is an identity matrix<sup>1</sup>.

<sup>&</sup>lt;sup>1</sup>In the implementation, we omit W instead.

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

157

158

159

160

161

162

163

164

166

167

168

169

170

171

172

173

174

The mapping function of student and teacher model is based on the best ablation results of Jiao et al. (2020), which defined as follows:

$$k = i \cdot \frac{N_T}{N_S}, \quad i \in [1, N_S], \quad k \in [1, N_T]$$
 (5)

In this formula, *i* represents the index of the student's layer, while k is the index of the corresponding teacher's layer.  $N_S$  denotes the total number of layers in the student model, and  $N_T$  represents the total number of layers in the teacher model.

The KD loss can be formulated as follows:

$$\mathcal{L_{KD}} = \mathcal{L}_{att} + \mathcal{L}_{hid} + \mathcal{L}_{embd} + \mathcal{L}_{pred}$$

We calculate the classification loss  $\mathcal{L}_{clf}$  as follows:

$$\mathcal{L}_{clf} = CE(z^S, GT)$$

Where GT is the ground truth of the observed instance.

Finally, we obtain the overall loss  $\mathcal{L}_{overall}$  which is going to be minimized in the training process:

$$\mathcal{L}_{overall} = \mathcal{L}_{\mathcal{KD}} + \mathcal{L}_{clf}$$

We use Mean Square Error (MSE) instead of KL Divergence due to faster convergence and higher performance, as supported by the experiment of Nityasya et al., 2022.

#### 2.2 Model Initialization

We use model initialization approach from DistilBERT (Sanh et al., 2020), where the student model's weights are initialized by copying the weights of the teacher model.

We alternately copy the weights of the teacher's embedding layer and classification layers to the student model. For the self-attention layer, we copy the weights based on the following mapping function:

$$\mathbf{SA}_T^j = \mathbf{SA}_S^{i*2} \quad \text{for} \quad i, j \in \mathbb{Z}^+$$
 (6)

Here, SA denotes the self-attention layers of the teacher T and student S models, respectively. The notations i and j indicate the indices of the student and teacher self-attention layers, respectively. To illustrate, the second self-attention layer of the teacher model will be mapped to the first self-attention layer of the student model. 175

If the student's hidden unit dimension is smaller than the teacher's, previous copy-weight approach will not work. Thus, we follow the approach of Xu et al., 2024, which has demonstrated a good initialization method for smaller model from large model by selecting evenly spaced elements in the teacher's linear weight and bias for the self-attention layer to map the student's self-attention layer correspondingly. For instance, suppose the teacher has a linear weight of 4x4, and the student has a 2x2 matrix; we select the 1st and 3rd slices along both the first and second dimensions. For the bias, we do the slicing in one dimension instead.

176

177

178

179

180

181

182

183

184

185

186

187

188

189

190

191

192

193

195

196

197

198

199

200

201

202

203

204

205

206

207

208

209

210

211

212

213

214

215

216

217

218

219

220

221

222

#### **Experiment Setup** 3

We provide three experiment setups: data, model, and training, that will be consistently used throughout this work.

**Data** We utilized massive (FitzGerald et al., 2022), Tweet Sentiment Multilingual dataset (denoted as tsm) (Barbieri et al., 2022), and universal-ner (Mayhew et al., 2024). We selected these datasets to observe the behavior of multilingual performance under different situations: high-resource data with parallel data (massive) and low-resource data with non-parallel data (tsm). We also use universal-ner for more complex multilingual task which has medium-resource data with non-parallel data. In this experiment, these data comprises of languages that are divided into unseen lang and seen lang to simulate zero-shot cross-lingual scenarios. Table 11 in Appendix B shows the corresponding datasets' data statistics and language partition. The detailed language partition information can be seen in Appendix C.

Model We used transformers library (Wolf et al., 2020) and the off-the-shelf implementation of xlm-roberta (Liu et al., 2019) and mdeberta (He et al., 2021) models. We used a reduction factor of 2 for the number of student layers compared to the teacher. Additionally, we compared performance by reducing the hidden units by half and keeping the hidden units the same as the teacher's. We experimented with three different model initialization scenarios: copying the weights from xlm-roberta-base or mdeberta-v3-base (from-base), copying from the fine-tuned teacher (from-teacher), and initializing without copying from any model (from-scratch). we compared their performances to understand the differences

Model	Method	Layers	XLM-R				Avg		
			massive	tsm	universal-ner	massive	tsm	universal-ner	8
Teacher	from-base	12	80.18	70.10	87.66	81.36	66.96	88.81	79.18
Student .	from-scratch from-scratch+KD	6 6	75.19 79.23	50.20 54.13	45.68 48.29	75.93 76.22	49.86 51.80	46.21 39.55	57.18 60.30
	from-teacher from-teacher + KD	6 6	81.18 <b>81.63</b>	62.99 <b>67.61</b>	79.50 <b>80.18</b>	80.45 <b>82.31</b>	58.12 61.08	78.20 <b>79.17</b>	73.41 <b>75.33</b>

Table 1: Comparative performance (F1-scores %) of XLM-R and mDeBERTa models across different datasets (all languages), initialization methods, and with/without knowledge distillation (KD). The teacher model has 12 layers, while all student models have 6 layers. All models are trained with all languages for each data. Bold denotes best methods for each data and their average.

between these strategies.

227

228

237

239

240

241

242

243

244

245

246

247

248

249

251

254

255

257

259

Training То fine-tune the model and perform knowledge distillation, we used AdamW (Loshchilov and Hutter, 2019) as the optimizer, with the default hyperparameters stated in the transformers library. We set the number of epochs to 30 and obtained the best results evaluated on the development set using the F1 score metric. The evaluation steps for massive are as follows: 5000 steps for §4.1 and §4.2, and 100 steps for §4.3. For tsm and universal-ner, we set the evaluation steps to 500, 250, and 60 for §4.1, §4.2, and §4.3, respectively. These differences are due to the data sizes used in the corresponding experiments. The rest of the hyperparameters follow the default configuration in the transformers library. We used an A100 GPU to train our models to run our experiments, running each model's training three times with different seeds. Performances are evaluated using Macro F1-Score. For universal-ner, it is evaluated using overall-f1 of seqeval wrapped in evaluate package<sup>2</sup>.

#### 4 Multilingual Transferability in KD

The ability of knowledge distillation (KD) to transfer knowledge across multiple languages efficiently remains unexplored. As mentioned in §2, two components need to be analyzed: model initialization and the distillation process itself. It is still unclear which of these factors contributes the most to the overall performance. Also, in multilingual scenarios, we often encounter situations where not all languages are covered in the training set. Understanding whether KD can facilitate zero-shot crosslingual (ZSCL) generalization and effectively trans-

<sup>2</sup>https://huggingface.co/spaces/ evaluate-metric/seqeval fer multilingual knowledge remains unexplored. Building a multilingual dataset is tedious, thus leading researchers to opt for using one language. As a result, it is desirable if such setup can achieve cross-lingual generalization (Artetxe and Schwenk, 2019).

261

262

263

265

266

267

268

269

270

271

272

273

274

275

276

277

278

279

281

282

283

285

287

288

290

291

292

293

294

295

# 4.1 Weight Copy Transfers More Information vs Distillation Loss

Table 1 presents a comparative analysis of model performance with and without the copy-weight initialization strategy using all languages in the training dataset. The results demonstrate that the from-teacher approach significantly outperforms the from-scratch method, particularly when trained on the tsm and universal-ner datasets. While the performance gap is less pronounced for the massive dataset, it still favors the from-teacher method. These findings suggest that directly copying weights from a larger model can serve as an effective initialization strategy, especially for low-resource scenarios.

Furthermore, the application of knowledge distillation consistently yields performance improvements across all configurations. However, it is noteworthy that the initial weight initialization plays a more critical role in determining overall model performance. The from-scratch initialization, even with knowledge distillation, struggles to match the performance levels achieved by the from-teacher.

To evaluate the effectiveness of different initialization strategies, we compared models initialized from the teacher to those initialized from the base model. As shown in Table 2, the from-teacher approach demonstrates a slight performance advantage over from-base, which can be attributed to the teacher model's prior fine-tuning. The application of Knowledge Distillation (KD) improves

Initialization		XLN	M-R		ERTa	Avg	
	massive	massive tsm universal-ner massive		tsm	universal-ner	.9	
from-base	80.37	60.17	78.64	80.61	58.24	77.04	72.51
+KD	81.61	65.94	<b>80.29</b>	81.82	59.11	<b>79.24</b>	74.67
from-teacher	81.18	62.99	79.50	80.45	58.29	78.20	73.44
+KD	<b>81.63</b>	<b>67.61</b>	80.18	<b>82.31</b>	<b>61.08</b>	79.17	<b>75.33</b>

Table 2: Performance comparison of XLMR and MDEBERTA models on from-base and from-teacher initialization. Bold denotes best methods for each data and their average.

Model	Hidden Size	XLM-R				Avg		
		massive	tsm	universal-ner	massive	tsm	universal-ner	
from-scratch	384	78.05	51.63	36.46	75.93	49.97	44.33	56.09
	768	79.23	54.13	48.29	76.22	49.86	52.13	59.98
from-teacher	384	79.90	58.34	44.87	80.10	54.75	44.87	60.47
	768	81.63	67.61	80.18	82.31	61.08	79.17	75.33

Table 3: Comparative performance of XLM-R and mDeBERTa models across different datasets, initialization methods, and hidden sizes using Knowledge Distillation. from-scratch and from-teacher use a layer reduction factor of 2.

performance for both initialization methods. For datasets such as massive and universal-ner, the performance gains are comparable between the two initialization strategies when KD is applied.

299

300

301

302

305

306

307

308 309

310

311

312

314

315

316

317

319

320

321

323

324

326

Our previous experiments simply reduced the model size by reducing the layer size while retaining the unit size. In practice, however, we might also want to reduce the unit size. In from-teacher initialization, the performance in all datasets falls significantly in tsm and universal-ner by halving the unit size. This performance reduction is expected due to the model's decreased capacity to retain multilingual information. Additionally, the current approach of uniformly copying weights (Xu et al., 2024) may not be optimal for this multilingual distillation task.

# 4.2 Knowledge of Unseen Languages is Transferrable with Seen Language Teacher Weight Copy

To test the zero-shot cross-lingual performance, we observe two conditions: 1) using the seen lang subset as training data for both the student and teacher models, and 2) using the seen lang subset as training data for the teacher, then using the English language to fine-tune the student model. The motivation is to observe if the model retains multilingual information from the copy-weight, even when fine-tuned using only English.

Table 4 shows the results of zero-shot crosslingual generalization. The teacher's accuracy drops significantly compared to the scenario in §4.1. When using the seen lang subset for training, we observe similar behavior to the previous results, with a slight difference between in with and without KD in both datasets, unlike the results shown in §4.1. However, without weight-copy, the performance of each data plummets to near-random answers. **This shows that weight-copy preserves multilingual knowledge and enables zero-shot cross-lingual generalization in both high and low-resource scenarios**. 328

329

330

331

332

333

334

335

336

337

338

339

340

341

342

343

344

345

347

349

350

351

352

353

354

355

356

357

Using English as the training data deteriorates the performance in massive's without KD setup, yet it gains considerable performance by using KD in the copy-weight initialization. Even when using only the English language, the student still retains the ZSCL generalization performance, albeit with reduced effectiveness, which further strengthens our claim regarding copy-weight multilingual generalization.

On the other hand, tsm and universal-ner perform similarly to the student model trained on seen lang, where using KD yields only a marginal increase. This is attributed to KD needing a sufficient amount of data to be effective.

# 4.3 Multilingual Distillation is Possible Even if Only English Data is Available

We fine-tune the teacher and student models using only English, as this language is the most widely available. Note that, unlike the experiment done in

Model	Initialization	Train Data	XLM-R				Avg		
			massive	tsm	universal-ner	massive	tsm	universal-ner	. 8
Teacher	from-base	seen-lang	68.22	57.11	84.38	68.54	57.04	87.50	70.47
Student	from-teacher from-scratch	seen-lang seen-lang	65.74 15.10	54.02 40.27	78.47 14.61	55.62 16.36	47.35 32.74	72.77 14.61	62.33 22.28
	from-teacher from-scratch	english english	59.65 7.55	53.35 33.24	73.73 7.09	32.82 3.92	46.18 28.89	66.61 7.50	55.39 14.70

Table 4: Cross-lingual performance (F1-scores %) for XLM-R and mDeBERTa models with different initialization and training data. Knowledge Distillation is used on these models. **Teacher models are trained using seen-lang**.

Model	Initialization	massive	tsm	universal-ner	Avg
XLMR (T)	from-base	56.42	58.97	77.86	64.42
XLMR (S)	from-teacher	47.85	54.18	69.05	57.03
XLMR (S)	from-scratch	7.50	35.10	9.23	17.28
mDEBERTa (T)	from-base	63.19	59.29	77.93	66.80
mDEBERTa (S)	from-teacher	26.83	44.54	61.97	44.45
mDEBERTa (S)	from-scratch	3.72	28.45	10.12	14.10

Table 5: Cross-lingual performance (F1-scores). (T) **denotes Teacher model trained on English only**; (S) denotes Student model.

§4.2, we do not train the teacher model with seen lang; we use English instead. We then evaluate it on unseen lang to make the results comparable with those in §4.2. We focus on KD since it has shown a consistent pattern in the previous experiments in §4.1 and §4.2.

Table 5 shows the results of the current experiment. Compared to using a fine-tuned teacher model with seen lang, we can see that massive performance dropped by about 12%, while it slightly improved for tsm and universal-ner. We argue that tsm and universal-ner data is nonparallel and contains substantial fewer instances than massive. As a result, these performances do not follow the same pattern as massive.

Although from-teacher exhibited the highest score, there is a significant gap compared to the §4.2 experiment. Having one language trained on the teacher makes copy-weight initialization less effective, yet the model still retains some multilingual capability. In contrast, from-scratch performs similarly and near random score<sup>3</sup>.

Training Method	massive	tsm	universal-ner
With Finetune	81.63	67.61	80.18
Without Finetune	38.05	33.57	7.79
Random Score	7.02	33.33	1.75

Table 6: Zero-shot performance by only copying the weight of the respective fine-tuned teacher to their half-layer students. Random score obtained on universal-ner is generated through average of random predictions in 30 runs.

# 5 Behavior Analysis in Copy-weight Strategy

In §4, we summarized that model initialization strategy significantly impacts transferring multilingual knowledge, with from-teacher performing the best. This section provides more detailed analysis related to the model's characteristics when using the copy-weight strategy, such as zero-shot copy classification performance (§5.1), training speed after the weight is copied from the teacher to the student model across different data subsets (§5.2), and 3) performance across different data subsets (§5.3).

The experiments performed in this section will use KD and the setup described in §4.1, with full hidden size. We focus on analyzing the behavior of the copy-weight strategy.

# 5.1 Weight Copy model preserve some information even without finetuning

Given the that copy-weight approach is better than the distillation technique itself, we investigate how much multilingual information is retained simply by copying the weights without any additional fine-tuning. Table 6 provides the performance results. We observe that these scores are substantially lower than when fine-tuning is performed. **Intriguingly, massive and universal-ner scores are not as low as random guesses, implying that** 

380 381

382

383

384

386

387

388

389

390

391

392

393

394

395

396

397

398

399

400

401

402

403

404

405

406

407

<sup>&</sup>lt;sup>3</sup>Random score is obtained by making all predictions equal to the major class in the dataset.



Figure 2: Performance across different data subsets in different initialization strategies.

some knowledge is still retained, though not fully 'connected,' and needs to be fine-tuned. On the other hand, tsm shows performance comparable to random guessing. We hypothesize that this is due to the low number of instances in tsm, which do not preserve the inherent bias of multilingual knowledge as strongly as the others.

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

429

430

431

432

434

435

437

439

441

#### 5.2 Weight Copy Models Achieve Higher **Performance with Less Data**

The experiment in \$4.1 demonstrated that the copyweight approach exhibited higher performance, especially in the massive dataset due to its large number of instances. We argue that these results are attributed to better initialization, which enhances data efficiency. To test this hypothesis, we conducted an experiment by creating four subsets of the massive dataset, consisting of 1%, 5%, 10%, and 20% of the original data. These subsets were generated using stratified sampling based on the label distribution for each language.

Figure 2 illustrates the results for the three 428 model initialization strategies. We observe a pattern where using more data corresponds to higher scores. The performance order is consistent, with the best scores achieved by from-teacher and the worst by from-scratch. In the 1% data sub-433 set, from-teacher achieved around 69% f1-score, showing a significant gap compared to the others, with more than a 20% difference. However, as the 436 dataset size increases, the gap between scenarios becomes smaller, yet from-teacher consistently 438 exhibits the best results. This demonstrates that utilizing the teacher's fine-tuned weights, even in 440 a low-resource setting, benefits from the inherited information, providing better scores. 442



Figure 3: Training loss plot per step across different data subsets.

#### Weight Copy Provides Better 5.3 **Initialization – Model Converged Faster**

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

So far, we have explored that different data subsets exhibit different performances across model initialization strategies. This might also correlate with the training speed due to a better start. Thus, we are also interested in exploring learning efficiency by comparing the learning speed of each strategy with each data subset.

The resulting score correlates with the learning speed, depicted in Figure 3. Using the full data subset, we observe that the order of scenarios sorted by learning speed is similar to the order in Figure 2. With smaller data subsets, the gap in training loss between each model is wide, with from-teacher showing the fastest convergence rate. As more data is added, the gap becomes smaller. We posit that this is attributed from the fine-tuned teacher weightcopy, making the model learns faster. Furthermore, this also indicates that copying the teacher's weights in low-resource settings not only improves the score but also accelerates learning speed, reducing the cost of training the model.

#### **Related Work** 6

Model Initialization Model initialization is crucial when training a model. Glorot and Bengio (2010) introduced a method to properly initialize the weights of neural networks using a normal distribution to avoid issues related to vanishing and exploding gradients during training. This approach has been extended by several others, such as He et al. (2015), Mishkin and Matas (2016), and Saxe

571

572

573

525

526

et al. (2014), to add robustness to gradient problems. While these methods address numerical instability, they do not incorporate inherent initial knowledge. Transfer learning (Zhuang et al., 2020; Howard and Ruder, 2018) provides a way to start training a model with better initialization with prior knowledge. We pre-train the model on unlabeled data and then fine-tune it on the desired task. Multilingual models like DeBERTa (He et al., 2021), mBERT (Devlin et al., 2019), and XLM-R (Liu et al., 2019) can be used to train models that handle multiple languages. However, these models require extensive training resources to create.

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

501

502

507

508

510

511

512

513

514

515

516

517

Knowledge Distillation Knowledge distillation (KD) (Hinton et al., 2015) produces models with fewer parameters (student model) guided by a larger model (teacher model), often resulting in higher quality than models trained from scratch. In NLP, KD can be applied directly to task-specific or downstream tasks (Nityasya et al., 2022; Adhikari et al., 2020; Liu et al., 2020b), or during the pre-training phase of the student model (Sun et al., 2020), which can then be fine-tuned. Several works apply KD during both pre-training and finetuning steps (Jiao et al., 2020; Sanh et al., 2019; Liu et al., 2020a). The aspects of the teacher model that the student should mimic can vary; a common approach is for the student to mimic only the probability distribution of the teacher's prediction layer output. However, Jiao et al. (2020) and Sun et al. (2020) also include outputs such as the teacher model's layer outputs, attention layers, and embedding layers. Wang et al. (2020) and Ansell et al. (2023) explore the potential of KD in multilingual settings, with the latter utilizing sparse fine-tuning and a setup similar to Jiao et al. (2020). However, these works do not thoroughly investigate the behavior of the approach, such as the impact of initialization and data size. This research work dives into probing the influence of these components, providing better insight into multilingual performance capability using weight-copy techniques.

### 7 Conclusion

518In this work, we observed the effectiveness of519Knowledge Distillation, in multilingual settings,520focusing on identifying the factors that signifi-521cantly influence the performance of student models.522Our findings demonstrated that model initializa-523tion, specifically through weight copying from a524fine-tuned teacher model, plays a crucial role in

enhancing the performance and learning speed of student models. This finding was consistent across both high-resource and low-resource datasets, highlighting the importance of weight initialization in retaining multilingual knowledge and facilitating effective KD.

These insights underscore the critical role of initialization in KD, suggesting that simple yet effective strategies, such as weight copying, can lead to substantial performance gains without requiring extensive data or computational resources. This work contributes valuable and practical insights to developing efficient and high-performing multilingual models, particularly in resource-constrained environments. A promising future work is to propose a novel, efficient initialization method that do selective weight-copy or pruning to have a better initialization for the distillation process or normal fine-tuning.

# 8 Limitations

In this work, we focus on sequence classification and token classification tasks, which may not generalize to other tasks, such as Natural Language Generation. The languages observed in this work are those represented in massive, universal-ner and tsm, which do not include every possible language. Additionally, our study is limited to specific model sizes and architectures (e.g., XLM-RoBERTa and mDEBERTa).

We concentrate on analyzing multilingual capability, specifically zero-shot fine-tuning generalization, rather than zero-shot inference as exhibited in Large Language Models. This focus is due to the amount of computational cost associated with fine-tuning such large models. Additionally, most tasks done in real case are mostly classification, which is the case why we focus on encoder models that are more suitable for these tasks.

Finally, while using unlabeled datasets for distillation may improve the system's performance, it adds another layer of complexity to our work. Analyzing the data for use in a multilingual setting is beyond the scope of this study. We leave this for future work.

# **Ethical Considerations**

This work has no ethical issues, as it focuses on analyzing the inner workings of a multilingual model in knowledge distillation. All artifacts used in this research are permitted for research purposes and

628

629

630

align with their intended usage in multilingualism.
Additionally, the data utilized does not contain
any personally identifiable information or offensive
content. We use AI Assistants (LLM and Grammarly) to assist our writing in correcting grammatical errors.

## References

581

582

583

584

585

586

588

589

590

591

592

594

595

596

597

599

602

604

610

611

612

614

615

616

617

618 619

621

623

625

626

- Ashutosh Adhikari, Achyudh Ram, Raphael Tang, William L. Hamilton, and Jimmy Lin. 2020. Exploring the limits of simple learners in knowledge distillation for document classification with DocBERT. In *Proceedings of the 5th Workshop on Representation Learning for NLP*, pages 72–77, Online. Association for Computational Linguistics.
  - Alan Ansell, Edoardo Maria Ponti, Anna Korhonen, and Ivan Vulić. 2023. Distilling efficient languagespecific models for cross-lingual transfer. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 8147–8165, Toronto, Canada. Association for Computational Linguistics.
- Mikel Artetxe and Holger Schwenk. 2019. Massively multilingual sentence embeddings for zeroshot cross-lingual transfer and beyond. *Transactions* of the Association for Computational Linguistics, 7:597–610.
- Francesco Barbieri, Luis Espinosa Anke, and Jose Camacho-Collados. 2022. XLM-T: Multilingual language models in Twitter for sentiment analysis and beyond. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 258–266, Marseille, France. European Language Resources Association.
- Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2020. Unsupervised Cross-lingual Representation Learning at Scale.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In 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), pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Jack FitzGerald, Christopher Hench, Charith Peris, Scott Mackie, Kay Rottmann, Ana Sanchez, Aaron Nash, Liam Urbach, Vishesh Kakarala, Richa Singh, Swetha Ranganath, Laurie Crist, Misha Britan, Wouter Leeuwis, Gokhan Tur, and Prem Natarajan. 2022. Massive: A 1m-example multilingual natural language understanding dataset with 51 typologically-diverse languages.

Xavier Glorot and Yoshua Bengio. 2010. Understanding the difficulty of training deep feedforward neural networks. In *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics* (*AISTATS*), volume 9, pages 249–256. JMLR Workshop and Conference Proceedings.

Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny Livshits, Danny Wyatt, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Francisco Guzmán, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Govind Thattai, Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan Misra, Ivan Evtimov, Jack Zhang, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Alwala, Karthik Prasad, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini, Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Kushal Lakhotia, Lauren Rantala-Yeary, Laurens van der Maaten, Lawrence Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo, Lukas Blecher, Lukas Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Mannat Singh, Manohar Paluri, Marcin Kardas, Maria Tsimpoukelli, Mathew Oldham, Mathieu Rita, Maya Pavlova, Melanie Kambadur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal, Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Ning Zhang, Olivier Duchenne, Onur Çelebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal Bhargava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic, Roberta Raileanu, Rohan Maheswari, Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ronnie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana Chennabasappa, Sanjay Singh, Sean Bell, Seohyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sha-

ran Narang, Sharath Raparthy, Sheng Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon Vandenhende, Soumya Batra, Spencer Whitman, Sten Sootla, Stephane Collot, Suchin Gururangan, Sydney Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas Scialom, Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh Ramanathan, Viktor Kerkez, Vincent Gonguet, Virginie Do, Vish Vogeti, Vítor Albiero, Vladan Petrovic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whitney Meers, Xavier Martinet, Xiaodong Wang, Xiaofang Wang, Xiaoqing Ellen Tan, Xide Xia, Xinfeng Xie, Xuchao Jia, Xuewei Wang, Yaelle Goldschlag, Yashesh Gaur, Yasmine Babaei, Yi Wen, Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao, Zacharie Delpierre Coudert, Zheng Yan, Zhengxing Chen, Zoe Papakipos, Aaditya Singh, Aayushi Srivastava, Abha Jain, Adam Kelsey, Adam Shajnfeld, Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand, Ajay Menon, Ajay Sharma, Alex Boesenberg, Alexei Baevski, Allie Feinstein, Amanda Kallet, Amit Sangani, Amos Teo, Anam Yunus, Andrei Lupu, Andres Alvarado, Andrew Caples, Andrew Gu, Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchandani, Annie Dong, Annie Franco, Anuj Goyal, Aparajita Saraf, Arkabandhu Chowdhury, Ashley Gabriel, Ashwin Bharambe, Assaf Eisenman, Azadeh Yazdan, Beau James, Ben Maurer, Benjamin Leonhardi, Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Hancock, Bram Wasti, Brandon Spence, Brani Stojkovic, Brian Gamido, Britt Montalvo, Carl Parker, Carly Burton, Catalina Mejia, Ce Liu, Changhan Wang, Changkyu Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai, Chris Tindal, Christoph Feichtenhofer, Cynthia Gao, Damon Civin, Dana Beaty, Daniel Kreymer, Daniel Li, David Adkins, David Xu, Davide Testuggine, Delia David, Devi Parikh, Diana Liskovich, Didem Foss, Dingkang Wang, Duc Le, Dustin Holland, Edward Dowling, Eissa Jamil, Elaine Montgomery, Eleonora Presani, Emily Hahn, Emily Wood, Eric-Tuan Le, Erik Brinkman, Esteban Arcaute, Evan Dunbar, Evan Smothers, Fei Sun, Felix Kreuk, Feng Tian, Filippos Kokkinos, Firat Ozgenel, Francesco Caggioni, Frank Kanayet, Frank Seide, Gabriela Medina Florez, Gabriella Schwarz, Gada Badeer, Georgia Swee, Gil Halpern, Grant Herman, Grigory Sizov, Guangyi, Zhang, Guna Lakshminarayanan, Hakan Inan, Hamid Shojanazeri, Han Zou, Hannah Wang, Hanwen Zha, Haroun Habeeb, Harrison Rudolph, Helen Suk, Henry As-743 pegren, Hunter Goldman, Hongyuan Zhan, Ibrahim Damlaj, Igor Molybog, Igor Tufanov, Ilias Leontiadis, Irina-Elena Veliche, Itai Gat, Jake Weissman, James Geboski, James Kohli, Janice Lam, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe Cummings, Jon Carvill, Jon Shepard, Jonathan Mc-Phie, Jonathan Torres, Josh Ginsburg, Junjie Wang, Kai Wu, Kam Hou U, Karan Saxena, Kartikay Khandelwal, Katayoun Zand, Kathy Matosich, Kaushik

694

702

711

713

714

716

718

719

721

723

725

726

727

728

731

733

735

736

737 738

739

740

741

742

744

745

746

747

748

749

751

752

753

Veeraraghavan, Kelly Michelena, Keqian Li, Kiran Jagadeesh, Kun Huang, Kunal Chawla, Kyle Huang, Lailin Chen, Lakshya Garg, Lavender A, Leandro Silva, Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madian Khabsa, Manav Avalani, Manish Bhatt, Martynas Mankus, Matan Hasson, Matthew Lennie, Matthias Reso, Maxim Groshev, Maxim Naumov, Maya Lathi, Meghan Keneally, Miao Liu, Michael L. Seltzer, Michal Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark, Mike Macey, Mike Wang, Miquel Jubert Hermoso, Mo Metanat, Mohammad Rastegari, Munish Bansal, Nandhini Santhanam, Natascha Parks, Natasha White, Navyata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier, Nikhil Mehta, Nikolay Pavlovich Laptev, Ning Dong, Norman Cheng, Oleg Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem Kalinli, Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip Bontrager, Pierre Roux, Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani, Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu Nayani, Rahul Mitra, Rangaprabhu Parthasarathy, Raymond Li, Rebekkah Hogan, Robin Battey, Rocky Wang, Russ Howes, Ruty Rinott, Sachin Mehta, Sachin Siby, Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov, Satadru Pan, Saurabh Mahajan, Saurabh Verma, Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lindsay, Shaun Lindsay, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha, Shishir Patil, Shiva Shankar, Shuqiang Zhang, Shuqiang Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen Chen, Steve Kehoe, Steve Satterfield, Sudarshan Govindaprasad, Sumit Gupta, Summer Deng, Sungmin Cho, Sunny Virk, Suraj Subramanian, Sy Choudhury, Sydney Goldman, Tal Remez, Tamar Glaser, Tamara Best, Thilo Koehler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim Matthews, Timothy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai Mohan, Vinay Satish Kumar, Vishal Mangla, Vlad Ionescu, Vlad Poenaru, Vlad Tiberiu Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Constable, Xiaocheng Tang, Xiaojian Wu, Xiaolan Wang, Xilun Wu, Xinbo Gao, Yaniv Kleinman, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi, Youngjin Nam, Yu, Wang, Yu Zhao, Yuchen Hao, Yundi Qian, Yunlu Li, Yuzi He, Zach Rait, Zachary DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang, Zhiwei Zhao, and Zhiyu Ma. 2024. The llama 3 herd of models.

754

755

756

758

761

762

763

764

765

766

769

771

772

774

775

776

779

781

782

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

- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2015. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. Proceedings of the IEEE international conference on computer vision (ICCV), pages 1026–1034.
- Pengcheng He, Xiaodong Liu, Jianfeng Gao, and Weizhu Chen. 2021. Deberta: Decoding-enhanced bert with disentangled attention.

920

921

922

923

924

925

926

927

928

929

930

931

872

817 818

820

825

826

827

828

829

831

832

835 836

837

838

839

840

841

842

843

850

851

852

855

856

857

866

867

871

816

- Geoffrey Hinton, Oriol Vinyals, and Jeffrey Dean. 2015. Distilling the knowledge in a neural network. In *NIPS Deep Learning and Representation Learning Workshop*.
- Jeremy Howard and Sebastian Ruder. 2018. Universal language model fine-tuning for text classification. In Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 328–339, Melbourne, Australia. Association for Computational Linguistics.
- Binyuan Hui, Jian Yang, Zeyu Cui, Jiaxi Yang, Dayiheng Liu, Lei Zhang, Tianyu Liu, Jiajun Zhang, Bowen Yu, Keming Lu, et al. 2024. Qwen2. 5-coder technical report. *arXiv preprint arXiv:2409.12186*.
- Xiaoqi Jiao, Yichun Yin, Lifeng Shang, Xin Jiang, Xiao Chen, Linlin Li, Fang Wang, and Qun Liu. 2020. TinyBERT: Distilling BERT for natural language understanding. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 4163– 4174, Online. Association for Computational Linguistics.
- Weijie Liu, Peng Zhou, Zhiruo Wang, Zhe Zhao, Haotang Deng, and Qi Ju. 2020a. FastBERT: a selfdistilling BERT with adaptive inference time. In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, pages 6035– 6044, Online. Association for Computational Linguistics.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019.
  Roberta: A robustly optimized bert pretraining approach. arXiv preprint arXiv:1907.11692.
- Yuang Liu, Wei Zhang, and Jun Wang. 2020b. Adaptive multi-teacher multi-level knowledge distillation. *Neurocomputing*, 415:106–113.
- Ilya Loshchilov and Frank Hutter. 2019. Decoupled weight decay regularization. In *International Conference on Learning Representations*.
- Stephen Mayhew, Terra Blevins, Shuheng Liu, Marek Šuppa, Hila Gonen, Joseph Marvin Imperial, Börje F. Karlsson, Peiqin Lin, Nikola Ljubešić, LJ Miranda, Barbara Plank, Arij Riab, and Yuval Pinter. 2024. Universal ner: A gold-standard multilingual named entity recognition benchmark. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics* (NAACL).
- Dmytro Mishkin and Jiri Matas. 2016. All you need is a good init. In *International Conference on Learning Representations (ICLR).*
- Made Nindyatama Nityasya, Haryo Akbarianto Wibowo, Rendi Chevi, Radityo Eko Prasojo, and Alham Fikri Aji. 2022. Which student is best? a comprehensive knowledge distillation exam for taskspecific bert models.

- Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. 2019. Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter. In *NeurIPSEMC*<sup>2</sup>Workshop.
- Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. 2020. Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter.
- Andrew M Saxe, James L McClelland, and Surya Ganguli. 2014. Exact solutions to the nonlinear dynamics of learning in deep linear neural networks. In *International Conference on Learning Representations* (*ICLR*).
- Zhiqing Sun, Hongkun Yu, Xiaodan Song, Renjie Liu, Yiming Yang, and Denny Zhou. 2020. MobileBERT: a compact task-agnostic BERT for resource-limited devices. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 2158–2170, Online. Association for Computational Linguistics.
- Gemma Team, Morgane Riviere, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhupatiraju, Léonard Hussenot, Thomas Mesnard, Bobak Shahriari, Alexandre Ramé, Johan Ferret, Peter Liu, Pouya Tafti, Abe Friesen, Michelle Casbon, Sabela Ramos, Ravin Kumar, Charline Le Lan, Sammy Jerome, Anton Tsitsulin, Nino Vieillard, Piotr Stanczyk, Sertan Girgin, Nikola Momchev, Matt Hoffman, Shantanu Thakoor, Jean-Bastien Grill, Behnam Neyshabur, Olivier Bachem, Alanna Walton, Aliaksei Severyn, Alicia Parrish, Aliya Ahmad, Allen Hutchison, Alvin Abdagic, Amanda Carl, Amy Shen, Andy Brock, Andy Coenen, Anthony Laforge, Antonia Paterson, Ben Bastian, Bilal Piot, Bo Wu, Brandon Royal, Charlie Chen, Chintu Kumar, Chris Perry, Chris Welty, Christopher A. Choquette-Choo, Danila Sinopalnikov, David Weinberger, Dimple Vijaykumar, Dominika Rogozińska, Dustin Herbison, Elisa Bandy, Emma Wang, Eric Noland, Erica Moreira, Evan Senter, Evgenii Eltyshev, Francesco Visin, Gabriel Rasskin, Gary Wei, Glenn Cameron, Gus Martins, Hadi Hashemi, Hanna Klimczak-Plucińska, Harleen Batra, Harsh Dhand, Ivan Nardini, Jacinda Mein, Jack Zhou, James Svensson, Jeff Stanway, Jetha Chan, Jin Peng Zhou, Joana Carrasqueira, Joana Iljazi, Jocelyn Becker, Joe Fernandez, Joost van Amersfoort, Josh Gordon, Josh Lipschultz, Josh Newlan, Ju yeong Ji, Kareem Mohamed, Kartikeya Badola, Kat Black, Katie Millican, Keelin McDonell, Kelvin Nguyen, Kiranbir Sodhia, Kish Greene, Lars Lowe Sjoesund, Lauren Usui, Laurent Sifre, Lena Heuermann, Leticia Lago, Lilly McNealus, Livio Baldini Soares, Logan Kilpatrick, Lucas Dixon, Luciano Martins, Machel Reid, Manvinder Singh, Mark Iverson, Martin Görner, Mat Velloso, Mateo Wirth, Matt Davidow, Matt Miller, Matthew Rahtz, Matthew Watson, Meg Risdal, Mehran Kazemi, Michael Moynihan, Ming Zhang, Minsuk Kahng, Minwoo Park, Mofi Rahman, Mohit Khatwani, Natalie Dao, Nenshad Bardoliwalla, Nesh Devanathan, Neta Dumai, Nilay

Chauhan, Oscar Wahltinez, Pankil Botarda, Parker 932 Barnes, Paul Barham, Paul Michel, Pengchong 933 Jin, Petko Georgiev, Phil Culliton, Pradeep Kup-935 pala, Ramona Comanescu, Ramona Merhej, Reena Jana, Reza Ardeshir Rokni, Rishabh Agarwal, Ryan Mullins, Samaneh Saadat, Sara Mc Carthy, Sarah Cogan, Sarah Perrin, Sébastien M. R. Arnold, Sebastian Krause, Shengyang Dai, Shruti Garg, Shruti 939 Sheth, Sue Ronstrom, Susan Chan, Timothy Jordan, 941 Ting Yu, Tom Eccles, Tom Hennigan, Tomas Kocisky, Tulsee Doshi, Vihan Jain, Vikas Yadav, Vilobh 942 Meshram, Vishal Dharmadhikari, Warren Barkley, 943 Wei Wei, Wenming Ye, Woohyun Han, Woosuk 945 Kwon, Xiang Xu, Zhe Shen, Zhitao Gong, Zichuan 946 Wei, Victor Cotruta, Phoebe Kirk, Anand Rao, Minh Giang, Ludovic Peran, Tris Warkentin, Eli Collins, Joelle Barral, Zoubin Ghahramani, Raia Hadsell, D. Sculley, Jeanine Banks, Anca Dragan, Slav Petrov, Oriol Vinyals, Jeff Dean, Demis Hassabis, Koray Kavukcuoglu, Clement Farabet, Elena Buchatskaya, 951 Sebastian Borgeaud, Noah Fiedel, Armand Joulin, 952 Kathleen Kenealy, Robert Dadashi, and Alek An-953 dreev. 2024. Gemma 2: Improving open language models at a practical size.

> Xinyu Wang, Yong Jiang, Nguyen Bach, Tao Wang, Fei Huang, and Kewei Tu. 2020. Structure-level knowledge distillation for multilingual sequence labeling. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 3317–3330, Online. Association for Computational Linguistics.

957 958

959

962

963

964

965

966

967

968

969

970 971

972

973

975

976

977

978 979

981

982

983

Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. 2020. Transformers: State-of-the-art natural language processing. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations, pages 38–45, Online. Association for Computational Linguistics.

- Zhiqiu Xu, Yanjie Chen, Kirill Vishniakov, Yida Yin, Zhiqiang Shen, Trevor Darrell, Lingjie Liu, and Zhuang Liu. 2024. Initializing models with larger ones. In *The Twelfth International Conference on Learning Representations*.
- Fuzhen Zhuang, Zhiyuan Qi, Keyu Duan, Dongbo Xi, Yongchun Zhu, Hengshu Zhu, Hui Xiong, and Qing He. 2020. A comprehensive survey on transfer learning.

Model	del Initialization		#U	mass	massive		tsm		universal-ner	
				NON	KD	NON	KD	NON	KD	
	from-base (T)	12	768	80.13	-	70.10	_	87.66	_	
	from-teacher	6	768	81.18	81.63	62.99	67.61	79.50	80.18	
VIMD	from-base	6	768	80.37	81.61	60.17	65.94	78.64	80.29	
ALW-K	from-scratch	6	768	75.19	79.23	50.20	54.13	45.68	48.29	
	from-teacher	6	384	78.20	79.90	56.74	58.34	47.15	44.87	
	from-base	6	384	77.55	79.41	55.03	58.15	45.94	44.24	
	from-scratch	6	384	75.27	78.05	50.01	51.63	40.80	36.46	
	from-base (T)	12	768	81.36	-	66.96	_	88.81	_	
	from-teacher	6	768	80.45	82.31	58.29	61.08	78.20	79.17	
mDoDEDTo	from-base	6	768	80.61	81.82	58.24	59.11	77.04	79.24	
IIIDEDERIA	from-scratch	6	768	75.93	76.22	49.86	51.80	46.21	44.56	
	from-teacher	6	384	79.76	80.10	54.59	54,75	62.57	59.72	
	from-base	6	384	78.49	78.35	53.94	46.20	62.81	59.24	
	from-scratch	6	384	76.22	75.93	49.97	43.39	45.55	44.33	

Table 7: F1-scores (%) for XLM-R and mDeBERTa models with different initializations, both with (KD) and without (NON) knowledge distillation, across three datasets. #L: Number of Layers, #U: Number of Units, T: Teacher.

Model	Initialization	Training Data	Test Data	massive		tsm		universal-ner	
				NON-KD	KD	NON-KD	KD	NON-KD	KD
	from-base(T)	seen lang	unseen lang	68.22	_	57.11	_	84.38	_
	from-teacher	seen lang	unseen lang	64.30 62.77	65.74	53.99	54.02	77.05	78.47
XLM-R	from-scratch	seen lang	unseen lang	15.08	04.82 15.10	37.93	40.27	15.58	14.61
	from-teacher from-base from-scratch	english english english	unseen lang unseen lang unseen lang	47.23 38.69 5.25	59.65 46.16 7.55	54.36 50.47 34.39	53.35 49.74 33.24	66.96 64.99 9.21	73.73 71.73 7.09
	from-base(T)	seen lang	unseen lang	68.54	_	57.04	_	87.50	_
mDeBERTa	from-teacher from-base from-scratch	seen lang seen lang seen lang	unseen lang unseen lang unseen lang	50.91 47.59 14.80	55.62 54.90 16.36	48.10 46.66 39.59	47.35 49.04 32.74	73.28 69.13 15.58	72.77 73.56 14.61
	from-teacher from-base from-scratch	english english english	unseen lang unseen lang unseen lang	20.98 16.23 7.09	32.82 33.64 3.92	43.14 43.38 38.45	46.18 48.32 28.89	59.10 53.37 8.26	66.61 65.34 7.50

Table 8: F1-scores (%) for zero-shot cross-lingual generalization in Knowledge Distillation for XLM-R and mDeBERTa models. We fine-tune the teacher model using seen lang and fine-tune the student model according to the training data provided. NON-KD follows the student model configuration. All models contain 6 layers, except for the teacher models which have 12 layers.

Code	Language	Code	Language	Code	Language
ceb	Cebuano	en	English	sk	Slovak
da	Danish	hr	Croatian	sr	Serbian
de	German	pt	Portuguese	SV	Swedish
ru	Russian	tl	Tagalog	zh	Chinese

Table 9: Languages in the universal-ner dataset

Code	Language	Code	Language	Code	Language
arabic english french	Arabic English French	german hindi italian	German Hindi Italian	portuguese spanish	Portuguese Spanish

T 1 1	10	т	•	. 1		1
Table	10.	Languages	1n	the	tsm	dataset
ruore	10.	Dunguuges	111	unc	COM	autuset

Information	massive	tsm	universal-ner
Number of Training Data Number of Classes Number of Languages unseen lang partition	11,514 60 52 "am-ET", "cy-GB", "af-ZA", "km-KH", "sw-KE", "mn-MN", "tl-PH", "kn-IN", "te-IN", "sq-AL", "ur-PK", "az-AZ", "ml-IN", "ms-MY", "ca-ES", "sl-SL", "sv-SE", "ta-IN", "nl-NL", "it-IT", "he-IL", "pl-PL", "da-DK", "nb-NO", "ro-RO", "th-TH", "fa-IR"	1,839 3 8 "arabic", "french", "hindi", "portuguese"	6,366* 7 9 "da-ddt", "pt-bosque", "sr-set", "sk-snk", "sv-talbanken"

Table 11: Data statistics for massive, tsm, and universal-ner. Each language consists of the same number of instances in both datasets, except universal-ner which number instance varies across languages. unseen lang denotes language subset used in the zero-shot cross-lingual experiment. The rest of the languages are categorized as seen lang. universal-ner data does not include partition that does not have training set. \* denotes the mean of number of instances across languages.

Code	Language	Code	Language	Code	Language
af-ZA	Afrikaans	it-IT	Italian	ru-RU	Russian
am-ET	Amharic	ja-JP	Japanese	sl-SL	Slovanian
ar-SA	Arabic	jv-ID	Javanese	sq-AL	Albanian
az-AZ	Azeri	ka-GE	Georgian	sv-SE	Swedish
bn-BD	Bengali	km-KH	Khmer	sw-KE	Swahili
ca-ES	Catalan	ko-KR	Korean	hi-IN	Hindi
zh-CN	Chinese (China)	lv-LV	Latvian	kn-IN	Kannada
zh-TW	Chinese (Taiwan)	mn-MN	Mongolian	ml-IN	Malayalam
da-DK	Danish	ms-MY	Malay	ta-IN	Tamil
de-DE	German	my-MM	Burmese	te-IN	Telugu
el-GR	Greek	nb-NO	Norwegian	th-TH	Thai
en-US	English	nl-NL	Dutch	tl-PH	Tagalog
es-ES	Spanish	pl-PL	Polish	tr-TR	Turkish
fa-IR	Farsi	pt-PT	Portuguese	ur-PK	Urdu
fi-FI	Finnish	ro-RO	Romanian	vi-VN	Vietnamese
fr-FR	French	he-IL	Hebrew	cy-GB	Welsh
hu-HU	Hungarian	hy-AM	Armenian	id-ID	Indonesian
is-IS	Icelandic	-			

Table 12: Language codes and corresponding language names in massive dataset.