# Dicta-LM 3.0: Advancing The Frontier of Hebrew Sovereign LLMs

**Shaltiel Shmidman[1,†], Avi Shmidman[1,2,†], Amir DN Cohen[3,‡] Moshe Koppel[1,2,†]**

[1]DICTA / Jerusalem, Israel
[2]Bar Ilan University / Ramat Gan, Israel
[3]OriginAI / Ramat Gan, Israel

[†]{shaltiel,avi,moishk}@dicta.org.il

[‡]amirc@originai.co

## Abstract

Open-weight LLMs have been released by frontier labs; however, sovereign Large Language Models (for languages other than English) remain low in supply yet high in demand. Training large language models (LLMs) for low-resource languages such as Hebrew poses unique challenges. In this paper, we introduce Dicta-LM 3.0: an open-weight collection of LLMs trained on substantially-sized corpora of Hebrew and English texts. The model is released in three sizes: 24B - adapted from the Mistral-Small-3.1 base model, 12B - adapted from the NVIDIA Nemotron Nano V2 model, and 1.7B - adapted from the Qwen3-1.7B base model. We are releasing multiple variants of each model, each with a native context length of 65k tokens; base model and chat model with tool-calling support. To rigorously evaluate our models, we introduce a new benchmark suite for evaluation of Hebrew chat-LLMs, covering a diverse set of tasks including Translation, Summarization, Winograd, Israeli Trivia, and Diacritization (*nikud*). Our work not only addresses the intricacies of training LLMs in low-resource languages but also proposes a framework that can be leveraged for adapting other LLMs to various non-English languages, contributing to the broader field of multilingual NLP.

## 1 Introduction

The development of generative language models has significantly advanced natural language processing (NLP), enhancing the sophistication and contextual understanding in human-computer interactions (OpenAI et al., 2024b; Team et al., 2024; Anthropic, 2024). Leading state-of-the-art open-weight generative LLM models, such as Qwen and Olmo (Yang et al., 2025; Groeneveld et al., 2024), are trained on trillions of tokens, but only a very small percentage of that data represents low-resource languages such as Hebrew (Touvron et al., 2023). Consequently, these models significantly underperform in such languages.

Training large language models (LLMs) in low-resource languages presents unique challenges, stemming from limited data availability, complex morphological structures, and the lack of robust evaluation frameworks tailored to these languages. Hebrew, with its rich morphology and limited large-scale corpora, exemplifies these difficulties.

In response to these challenges, we introduce Dicta-LM 3.0, a collection of generative language models specifically optimized for Hebrew. The models are trained with a native context length of 65k, and we release them in three sizes. The first is 24B parameters, built upon the Mistral-Small 3.1 model (Mistral AI, 2025). The second is 12B parameter, built upon NVIDIA Nemotron Nano V2 (NVIDIA et al., 2025). The third is a 1.7B parameter built upon Qwen3 1.7B (Yang et al., 2025). All three models were continuously pre-trained on approximately 100 billion tokens of Hebrew, mixed together with 30B tokens of English data.

For each model size, we release multiple variants of the model: the base model and a chat model with tool-calling capabilities. We evaluate the base models on the Hebrew LLM Leaderboard (Shmidman et al., 2024), a collection of Hebrew benchmarks evaluating base models using few-shot prompting (Brown et al., 2020). To address the evaluation gap for non-base Xomodels, we introduce a new benchmark suite for assessing chat-style Hebrew language models. This suite includes tasks such as Translation, Summarization, Wingorad, Israeli Trivia, and Diacritization (*nikud*). Our comprehensive evaluation demonstrates that DictaLM3.0 achieve state-of-the-art performance on these tasks for its weight class.

The methodologies and evaluation frameworks presented in this work provide insights and potential pathways for adapting other LLMs to various non-English languages. This research contributes to the

broader field of multilingual NLP by addressing the unique challenges posed by low-resource languages and offering scalable solutions for their integration into advanced language models.

## 2 Pre-training Data

Our pre-training corpus consisted of 75% Hebrew and 25% English data, detailed below.

### 2.1 Hebrew Data

#### 2.1.1 Sources

The Hebrew data consists of ~100B tokens, collected from a wide range of sources including available open-source corpora, internal web scraping, Hebrew books scanned and digitized in-house, and partnerships with companies that graciously provided their internal data for training. The data can be broken down into the following categories:
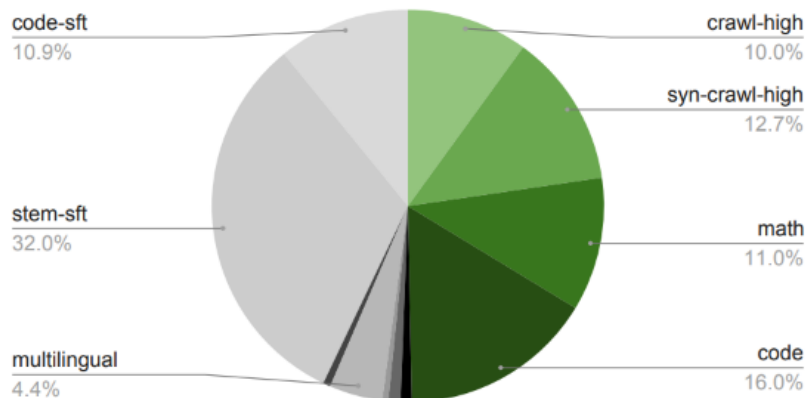
- **Internet**: Approximately 65% of the data. This category includes web crawls such as C4 (Raffel et al., 2019), OSCAR (Ortiz Suárez et al., 2019), FineWeb2 (Penedo et al., 2025), Wikipedia, and others. We processed over 150TB of unprocessed crawl data when extracting this data.

- **Social Media**: Approximately 17% of the data. This data includes various social media corpora such as Hebrew twitter posts and online Hebrew blogs.

- **News & Legal**: Approximately 9% of the data. This includes human transcriptions of TV and Radio, various news sites, Israeli parliament transcripts, and online legal resources.

- **Academia & Literature**: Approximately 8.5% of the data. This data includes academic papers, journals, books, and historical Hebrew texts (such as the Ben Yehuda project (Ben-Yehuda, 2024) and the Sefaria project (Sefaria, 2024)).

- **Tagged Data**: Approximately 0.5% of the data. This includes corpora tagged for various objectives such as NER, UD Dependencies (Tsarfaty, 2013; McDonald et al., 2013; Zeldes et al., 2022) and diacritics, as well as language resources such as dictionaries and thesauri. The data was reformatted into natural language explanatory text for the pretraining.

### 2.2 English Data

We include a 25% mix of English data in our pretraining corpus, in order to mitigate catastrophic forgetting and to enable the model to apply the logic and reasoning present in the English training corpora to Hebrew input texts. For our training, we picked out three corpora:

- **Nemotron-CC**: This is the pretraining corpus used to train the NVIDIA Nemotron Nano V2 model (NVIDIA et al., 2025). We sample a mixture from the various subsets of this corpus, using the same parameters used in Phase 3 of the Nano V2 base model training - displayed in Figure 1.

- **FineWeb-Edu** (Lozhkov et al., 2024): "1.3 trillion tokens of the finest educational data the web has to offer".

- **SlimPajama** (Soboleva et al., 2023): "extensively deduplicated, multi-corpora, open-source dataset for training large language models."

For the continued pre-training of the Nemotron Nano V2 model, we *only* used the Nemotron-CC data, whereas for the 24B model we split the English data 50/50 between Nemotron-CC and FineWeb-Edu. For the continued pre-training of the 1.7B model, we used the SlimPajama data mixed with 30% data from the Llama Nemotron Post Training Dataset (Bercovich et al., 2025) in straight text format (not chat template), as the Nemotron-CC data had not been released yet.

(c) Data mixture of Phase 3.

Figure 1: Data mixture used in Phase 3 of the Nemotron Nano V2 base model training. We used the same mixture in our CPT phase when mixing in the Nemotron-CC data.

## 3 Continuous Pre-training

Leading state-of-the-art generative LLM models are trained with trillions of tokens (Gemma Team et al., 2024; AI@Meta, 2024), making their reproduction very expensive. Therefore, we initialize our model from an existing SOTA model and continue pretraining it, leveraging the vast amount of data the model was already trained on. We chose to initialize from the following models, continuing to train each one separately:

- **Mistral-Small-3.1-24B-Base-2503**[1] (Mistral AI, 2025)

- **NVIDIA-Nemotron-Nano-12B-v2-Base**[2] (NVIDIA et al., 2025)

- **Qwen3-1.7B-Base**[3] (Yang et al., 2025)

Byte-pair encoding (BPE) tokenizers in foundation models often inadequately cover low-resource languages due to their scarcity in training data, leading to common words being split into many sub-tokens. This reduces compression, and therefore when choosing baseline models we had to ensure that the models had a good ratio, making them more practical for Hebrew.

We conducted the continuous pre-training on a compute cluster of 80 H200 141GB GPUs on NVIDIA DGX Cloud Lepton (Anand and Soman, 2025). All training was done using the NVIDIA NeMo Framework (Harper et al.), a scalable generative AI framework built for researchers and developers working on large language models, optimized for large-scale model training on NVIDIA hardware. Lepton's unified AI platform made it simple to scale distributed training across H200 GPUs and manage experiments through NeMo-Run, cutting iteration time and infrastructure overhead.

We split the continuous pre-training into two phases. In the first stage we iterated over our entire dataset, with a context window of 4,096 tokens. In the second stage we re-iterated over 20% of the dataset, with an extended context window of 65k tokens.

### 3.1 Phase 1

In the first phase we iterated over our entire dataset - 75% Hebrew and 25% English - summing to a total of 130B tokens. More details about the training are listed in Table 1; the training loss graph is shown in Figure 2.

---

[1] https://huggingface.co/mistralai/Mistral-Small-3.1-24B-Base-2503
[2] https://huggingface.co/nvidia/NVIDIA-Nemotron-Nano-12B-v2-Base
[3] https://huggingface.co/Qwen/Qwen3-1.7B-Base
[4] GBS differs because training capacity increased from 64 to 80 devices, requiring GBS to be divisible by 80.

| Hyperparameter | 24B | 12B | 1.7B |
|---|---|---|---|
| Global batch size (GBS)[4] | 256 | 320 | 256 |
| Micro batch size (MBS) | 4 | 1 | 8 |
| Tensor Parallel (TP) | 4 | 1 | 1 |
| Learning rate (LR) | 5e-6 | 5e-6 | 1e-5 |
| Sequence length | 4096 | | |
| Warmup steps | 2000 | | |
| Optimizer | Distributed AdamW | | |
| Schedule | Cosine annealing | | |
| Adam betas | $\beta_1 = 0.9, \ \beta_2 = 0.95$ | | |
| Clip grad | 1.0 | | |
| Weight decay | 0.1 | | |
| Total tokens | $\sim 130B$ | | |

Table 1: Training hyperparameters for phase 1 of the continuous pre-training.

## 3.2 Phase 2

In the second phase we extended the model's context length to 65k tokens. For this phase, we split our data into documents with >6144 tokens and <6144 tokens. We then continued to train the model with the extended context length, sampling 75% of the documents from the longer context pool, and 25% of the documents from the shorter context pool. We maintained the Hebrew-English ratio of 75/25, training for a total of 18B tokens. The details for this phase of the training are detailed in Table 2.

| Hyperparameter | 24B | 12B | 1.7B |
|---|---|---|---|
| Learning rate (LR) | 1e-6 | 1e-7 | 6e-7 |
| Tensor Parallel (TP) | 4 | 1 | 1 |
| Context Parallel (CP) | 5 | 16 | 8 |
| Micro batch size (MBS) | 1 | 1 | 2 |
| Global batch size (GBS) | 80 | 80 | 16 |
| Sequence length | 65,280 | 65,280 | 62,080 |
| Warmup steps | 200 | 200 | 1% |

Table 2: Training hyperparameters for phase 2 of the continuous-pretraining.

## 4   Base Model Evaluation

We evaluated the base models on Hebrew benchmarks from the Hebrew LLM-Leaderboard (hebrew-llm-leaderboard, 2025; Shmidman et al., 2024). We observe significant gains relative to the base models we initialized from, in the 24B, 12B, and 1.7B models. A full screenshot from the leaderboard including the 24B and 12B new base models is displayed in Figure 3. A before/after comparison of the three models can be found in Table 3. Notably, the 24B base model matches or exceeds the performance of models over four times larger, while also attaining the top score in the Israeli Trivia category, highlighting the effectiveness of sovereign LLMs.

We also evaluated the model on English benchmarks, in order to make sure the model retained the knowledge and abilities it was initially trained on. We evaluated the model on 3 popular benchmarks: Commonsense QA (Talmor et al., 2019), WinoGrande (Sakaguchi et al., 2019) and Arc-Challenge (Clark et al., 2018). The evaluation was done using the `lighteval` library (Habib et al., 2023). The exact command that we ran can be found in Appendix A. As can be seen in the table, we managed to retain >98% of the English capabilities of the original models.

Both base models are available for use on HuggingFace with a permissive license.
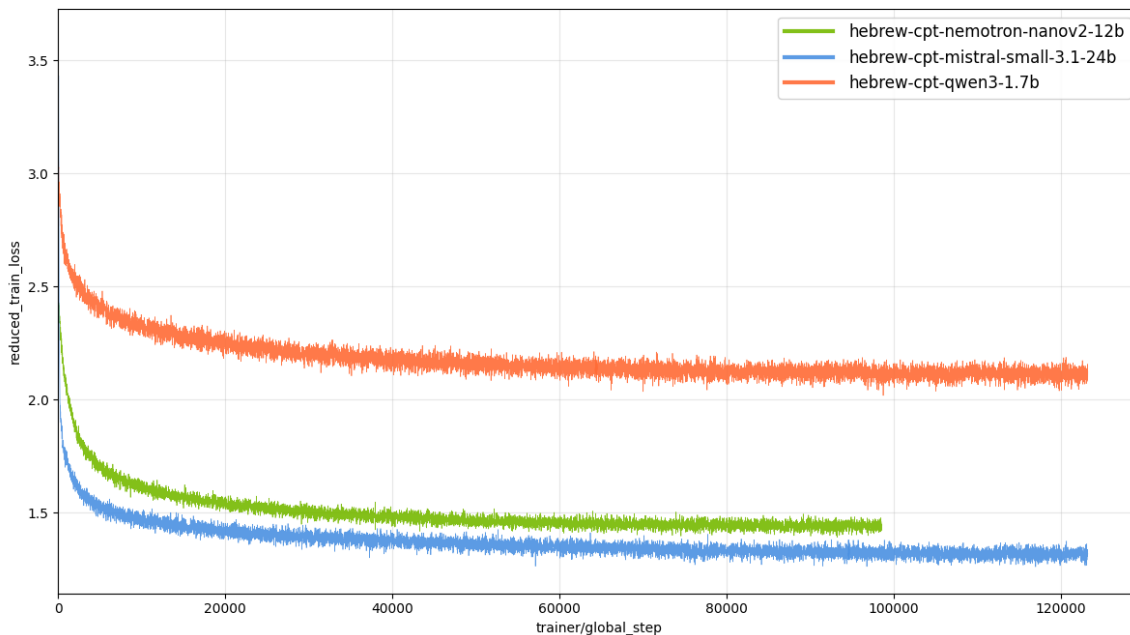
Figure 2: This graph depicts the loss value during the continuous pre-training stage; the 12B model had a larger batch size, which required fewer steps.

## 5 Post-training

Following the creation of our `DictaLM-3.0` base models via continuous pre-training, we produced variants of the models - chat models - by post-training the models on conversational data.

As part of our post-training pipeline, we produce a chat model - either instruct (GPT-3.5 style), or thinking (o1 style).

### 5.1 Supervised Fine-Tuning (SFT)

#### 5.1.1 SFT Data

For the first stage of post-training, we curated two conversational datasets for supervised fine-tuning. The first is the 'instruct' dataset, intended for regular chat conversations with direct answers. The second is the 'thinking' dataset, which adds a designated reasoning section before each assistant response.

For this stage, we aim for a 50/50 mix of English and Hebrew data. For most datasets, we obtain the Hebrew versions by translating them from their corresponding English versions with the `DictaLM-3.0` base model, using few-shot prompting. We also apply a filter to remove conversations that cannot be meaningfully translated (e.g., translation-request dialogues that would collapse into identical Hebrew text and introduce noise).

We gather our data from various open-source datasets. Following the findings by Shmidman et al. (2025), we use `gpt-oss-120b` (OpenAI, 2025) to regenerate the responses of the reasoning datasets. In addition, all datasets undergo post-processing to normalize formats and filter conversations through a rule-based system (e.g., removing refusals).

Below is an itemized list of the datasets in our final post-training SFT corpus. For datasets used for both thinking and instruct, the instruct version is derived by simply removing the reasoning section.

- Hermes 3 Dataset (NousResearch, 2024): We use 200k sampled conversations and include them along with their Hebrew translations. *Instruct only*.

- Hermes 3 - Augmented Subset: We sample 70k conversations that include a system message, take the first user prompt, and generate multi-turn dialogues using `gpt-oss-120b` (the model generating both stages separately). We then translate all conversations to Hebrew and interleave the languages, often adding small augmentations (e.g., adding 'Answer only in Hebrew' to the system prompt, adding

| Model | SNLI | QA | Sentiment | Winograd | Translation | IL-Facts | Average |
|---|---|---|---|---|---|---|---|
| *24B Parameter Models* | | | | | | | |
| Mistral-Small-3.1-24B | 81.2 | **77.9** | 66.7 | 80.9 | 30.8 | 58.5 | 66.0 |
| DictaLM-3.0-24B | **86.0** | 76.6 | **68.7** | **83.8** | **37.2** | **82.7** | **72.5** |
| *Improvement* | +4.8 | -1.3 | +2.0 | +2.9 | +6.4 | **+24.2** | **+6.5** |
| *12B Parameter Models* | | | | | | | |
| Nemotron-Nano-12B-v2 | 71.4 | 66.1 | 64.5 | 64.0 | 21.5 | 34.9 | 53.7 |
| DictaLM-3.0-Nemotron-12B | **80.5** | **77.8** | **72.4** | **79.5** | **33.9** | **54.8** | **66.5** |
| *Improvement* | +9.1 | +11.7 | +7.9 | +15.5 | +12.4 | **+19.9** | **+12.8** |
| *1.7B Parameter Models* | | | | | | | |
| Qwen3-1.7B-Base | 46 | 61.4 | 53 | 54.3 | 20.1 | 24.9 | 43.3 |
| DictaLM-3.0-1.7B | **67.1** | **62.8** | **61.9** | **57.2** | **29.2** | **30.6** | **51.5** |
| *Improvement* | +21.1 | +1.4 | +8.9 | +2.9 | +9.1 | +5.7 | **+8.2** |

Table 3: Performance comparison of the base models on Hebrew benchmarks before and after continuous pre-training

| Model | CommonSense | WinoGrande | Arc-Challenge | Avg. |
|---|---|---|---|---|
| *24B Parameter Models* | | | | |
| Mistral-Small-3.1-24B | 83.2 | 80.0 | 65.1 | 76.1 |
| DictaLM-3.0-24B | 82.8 | 78.2 | 63.4 | 74.8 |
| *Difference* | -0.4 | -1.8 | -1.7 | -1.3 |
| *12B Parameter Models* | | | | |
| Nemotron-Nano-12B-v2 | 79.4 | 75.7 | 64.7 | 73.3 |
| DictaLM-3.0-Nemotron-12B | **80.7** | **76.2** | 64.4 | **73.8** |
| *Difference* | +1.2 | +0.5 | -0.3 | +0.5 |
| *1.7B Parameter Models* | | | | |
| Qwen3-1.7B-Base | 74.4 | 61.5 | 38.8 | 58.2 |
| DictaLM-3.0-1.7B | 70.2 | 59.8 | **50.6** | **60.2** |
| *Difference* | -4.2 | -1.7 | +11.8 | +2.0 |

Table 4: Performance comparison of the base models on English benchmarks

'answer in English,' or simple context-switching to teach the model to follow the user's language). *Thinking only*.

- Math (Shmidman et al., 2025): We use 250k sampled conversations from the Nemotron Post Training Dataset (Nathawani et al., 2025; Bercovich et al., 2025), with a verified answered generated via `gpt-oss-120b`. We include them along with their Hebrew translations. *Instruct & Thinking*.

- rStarCoder (Liu et al., 2025): In this work we did not focus on code-generation abilities, since code is written in English and our emphasis is on Hebrew. Nevertheless, we include a small sample of code-related conversations, as previous work shows this can improve overall model capabilities. We use 30k sampled conversations and include them along with their Hebrew translations. *Instruct & Thinking*.

- Everyday & System conversations (SmolTalk v2 (Bakouch et al., 2025)): This dataset includes 30k examples of both everyday conversations and dialogues with a system prompt (teaching the model to adjust behavior based on the system message). The thinking subsets were regenerated using `gpt-oss-120b`. We include them along with their Hebrew translations. *Instruct & Thinking*.

- Tulu3 SFT Personas IF (Lambert et al., 2024): We do not include Hebrew translations for this dataset, as it contains very specific style constraints that do not translate well (e.g., requiring responses in all caps). *Instruct only*.

| T | Model | Average… | SNLI Ac… | QA TLNL… | Sentime… | Winogra… | Transla… | Israeli… | #Params… |
|---|---|---|---|---|---|---|---|---|---|
| ○ | mistralai/Mistral-Large-Instruct-2407 | 75.71 | 90.95 | 87.16 | 75.5 | 88.49 | 34.73 | 77.41 | 122.61 |
| ○ | ai21labs/AI21-Jamba-Large-1.6 | 73.42 | 93.1 | 81.5 | 68.42 | 81.65 | 35.98 | 80.07 | 398.56 |
| ● | meta-llama/Meta-Llama-3.1-70B | 72.57 | 88.57 | 80.31 | 77.5 | 80.94 | 34.7 | 73.42 | 70.55 |
| ● | dicta-il/DictaLM-3.0-24B-Base | 72.49 | 85.95 | 76.58 | 68.7 | 83.81 | 37.19 | 82.72 | 23.57 |
| ○ | CohereForAI/c4ai-command-r-plus-08-2024 | 72.36 | 91.67 | 82.45 | 67.83 | 85.25 | 36.85 | 70.1 | 103.81 |
| ○ | CohereForAI/c4ai-command-r-plus | 71.76 | 91.67 | 82.71 | 65.9 | 84.53 | 36 | 69.77 | 103.81 |
| ● | meta-llama/Meta-Llama-3-70B | 71.36 | 87.86 | 78.38 | 76.7 | 81.29 | 34.47 | 69.44 | 70.55 |
| ■ | meta-llama/Llama-3.3-70B-Instruct | 70.53 | 77.14 | 84.38 | 77.33 | 81.65 | 33.91 | 68.77 | 70.55 |
| ○ | CohereForAI/aya-expanse-32b | 70.24 | 92.14 | 83.42 | 67.97 | 82.73 | 34.7 | 60.47 | 32.3 |
| ● | Qwen/Qwen2.5-72B | 70.18 | 95.48 | 87 | 71.5 | 74.82 | 32.79 | 59.47 | 72.71 |
| ○ | nvidia/Llama-3.1-Nemotron-70B-Instruct-HF | 69.92 | 77.86 | 82.51 | 75.5 | 83.81 | 33.39 | 66.45 | 70.55 |
| ○ | meta-llama/Llama-3.1-70B-Instruct | 69.54 | 74.52 | 83.88 | 75.53 | 83.81 | 33.71 | 65.78 | 70.55 |
| ● | google/gemma-3-27b-pt | 69.5 | 85.24 | 78.27 | 64.97 | 81.65 | 36.45 | 70.43 | 27.43 |
| ● | meta-llama/Llama-4-Scout-17B-16E | 68.71 | 90 | 78.81 | 67.83 | 79.14 | 35.34 | 61.13 | 108.64 |
| ○ | CohereForAI/c4ai-command-r-08-2024 | 67.71 | 79.76 | 76.19 | 67 | 82.37 | 35.81 | 65.12 | 32.3 |
| ○ | mistralai/Mistral-Small-24B-Instruct-2501 | 67.56 | 83.1 | 65.97 | 71.67 | 79.5 | 29.95 | 55.15 | 23.57 |
| ● | dicta-il/DictaLM-3.0-Nemotron-12B-Base | 66.47 | 80.48 | 77.77 | 72.37 | 79.5 | 33.87 | 54.82 | 12.3 |
| ● | ibm-granite/granite-4.0-h-small-base | 66.39 | 88.81 | 74.27 | 71.67 | 75.9 | 34.18 | 53.49 | 32.21 |
| ■ | Qwen/QwQ-32B-Preview | 66.04 | 87.38 | 83.62 | 73.87 | 77.34 | 29.49 | 44.52 | 32.76 |
| ● | mistralai/Mistral-Small-3.1-24B-Base-2503 | 66 | 81.19 | 77.86 | 66.73 | 80.94 | 30.81 | 58.47 | 24.01 |

Figure 3: Up to date screenshot from the Hebrew LLM-Leaderboard, as of the publication date

- CCMatrix Translation (Schwenk et al., 2020): We sample 300k English–Hebrew pairs and convert them into conversational data using ~30 different templates. *Instruct only*.

- Agent-Ark/Toucan-1.5M (Xu et al., 2025): We sample 65k conversations from the SFT subset for the instruct data, and 50k conversations from the OSS subset for the thinking data. We include all conversations along with their Hebrew translations, with post-processing to ensure that all tool-calling JSON outputs are valid and match the expected schema. *Instruct & Thinking*.

- Identity: We auto-generated 3k conversations that include information about the model's identity (e.g., 'Who are you?' -> 'I am Dicta-LM..'). *Instruct & Thinking*.

- Others: We include about 100k other conversations that were synthesized directly by us. This includes mostly taking supervised datasets and converting them to conversational data using templates. For instance, given an annotated Hebrew sentence with morphological tagging, a conversation is generated in which a request is made regarding the morphological characteristics of a word with the sentence, and the response provides the relevant answer.

The final combined dataset for instruct consisted of ~2B tokens, from 1.5M conversations. The final combined dataset for reasoning consisted of ~3.2B tokens, from 725k conversations.

For the chat template, we drew inspiration from various existing models. For the message delimiter tokens we followed the convention of the Qwen3 models (Yang et al., 2025), for tool calling we followed the Hermes convention (Teknium et al., 2024), and for reasoning we followed the DeepSeek-R1 convention (DeepSeek-AI et al., 2025). We allocated the following special tokens in the tokenizer: `<|im_start|>`, `<|im_end|>`, `<tool_response>`, `</tool_response>`, `<tool_call>`, `</tool_call>`, `<think>`, `</think>`.

### 5.1.2 Training

We conducted the supervised fine-tuning on the same compute cluster on NVIDIA DGX Cloud Lepton, using NVIDIA NeMo Framework (Harper et al.).

We packed all the data into sequences of 65k tokens using the first-fit-decreasing method. For this phase of training, we computed the loss only on the completion tokens (the assistant responses). In addition, to maintain the knowledge gained during the pre-training phase, we mix in 10% pre-training data.

We first train the model for 10 epochs on the instruct data, and then we subsequently train the model for 5 epochs on the reasoning data. The hyperparameters can be found in Table 5.

| Hyperparameter | 24B | 12B | 1.7B |
|---|---|---|---|
| Learning rate (LR) | 5e-6 | 1e-6 | 1e-5 |
| Tensor Parallel (TP) | 4 | 1 | 1 |
| Context Parallel (CP) | 5 | 16 | 4 |
| Sequence length | 65,280 | 65,280 | 62,080 |
| Global batch size (GBS) | 40 | 48[5] | 48 |
| Micro batch size (MBS) | | 1 | |
| Warmup ratio | | 0.03 | |

Table 5: Training hyperparameters for the SFT stage.

## 5.2 Direct Preference Optimization (DPO)

The next stage in our model refinement process was Direct Preference Optimization (DPO) (Rafailov et al., 2023). DPO focuses on optimizing the model based on user preferences and feedback, enhancing its ability to generate responses that are both accurate and aligned with user expectations. We focused DPO on instruct variants due to the limited availability of diverse models employing the `gpt-oss` reasoning style, which is necessary for the comparative preference DPO requires.

### 5.2.1 DPO Data

We started off with the `llama-3.1-tulu-3-70b-preference-mixture` dataset (Lambert et al., 2024), which contains close to 300k preference pairs. We start by filtering out any pair that wasn't in English, and then sampling 150k examples from the dataset. As with the SFT data, we translated the 150k examples to Hebrew using the `Dicta-LM 3.0` 24B base model, so that our final dataset consists of a 50/50 mix of English and Hebrew.

As shown by (Lambert et al., 2024), it is important to mix on-policy examples as well. Therefore, we sampled 15% of our final dataset, and generated a response by passing the prompt to our 24B-Instruct model. We then used GPT-4o (OpenAI et al., 2024a) to compare our on-policy generated response with the original preferred response, to determine whether our response is preferred.

In addition, we mix in a few thousand examples of the following two types of behavioral preference data:

- We compiled the identity examples from our SFT corpus, and then used an LLM to generate negative examples (e.g., `You are ChatGPT, ...`) reinforcing the Dicta-LM identity in the model.

- We observed an interesting behavior during testing: after receiving a large tool response (such as a web-search result) in a language different from the original prompt, the model would continue responding in that new language. To address this, we generated synthetic preference pairs by sampling from the model, appending the instruction "Make sure your final response is in Hebrew" (or English) after the tool response, and then removing that instruction from the final preference pair.

As these two types of preference data were generated by artifically creating negative examples, we also generated similar preference pairs with the thinking section using our thinking model.

The final dataset consists of 260k examples.

### 5.2.2 Training

We conducted the supervised fine-tuning on the same compute cluster on NVIDIA DGX Cloud Lepton, using the NVIDIA **NeMo-RL** Training Framework (Nvidia, 2025). The NeMo-RL library integrated directly with Lepton, providing us with a seemless experience for scaling up post-training.

---

[5]Trained on 96 GPUs; global batch size chosen to remain divisible by 6.

We ran the DPO training using full parameter training (as opposed to PEFT), utilizing Megatron-LM as the backend (Shoeybi et al., 2020) achieving a near 2x speedup. For the instruct models, we run a full epoch on top of the SFT model. For the thinking models, we run 30 training steps on the two behavioral subsets on top of the GRPO model (see below). The full hyperparameters are listed in Table 6.

| Hyperparameter | Value |
|---|---|
| Max total sequence length | 4,096 |
| Global batch size (GBS) | 144 |
| Learning rate (LR) | 5e-7 |
| Minimum LR | 5e-8 |
| LR decay style | cosine |
| LR warmup iterations | 50 |
| Tensor Parallel (TP) | 4 |

Table 6: Training hyperparameters for DPO training.

## 5.3 Group Relative Policy Optimization (GRPO)

The next stage in our model refinement process was Group Relative Policy Optimization (GRPO) (Shao et al., 2024). GRPO directly trains the policy using grouped comparisons, improving the model's reasoning quality and overall response reliability. We applied GRPO to our thinking-style variants, as recent work has shown that GRPO often yields the largest improvements in reasoning and alignment quality.

### 5.3.1 GRPO Data

A key part of GRPO training is being able to calculate a reward for a generated response to a given prompt. In order to be able to accomplish that, we carefully curate a dataset, where we include metadata with each prompt allowing us to verify the correctness of a given answer and assign a numerical reward (between 0 and 1). Our dataset comprises of the following tasks:

- Instruction-Following (Zhou et al., 2023): This task defines a set of instructions which are appended to other prompts, defining how the response should look (e.g., no capital letters). Each prompt includes a verifier function which can be run on the final response to assess whether the constraint was met. We use the `allenai/RLVR-IFeval` dataset for English examples. For Hebrew, we manually curate a set of 30 Hebrew-specific constraints, and then sample prompts from our SFT data and randomly append one of our Hebrew-specific constraints to each prompt. The final score is either 0 (constraint wasn't followed) or 1 (constraint was followed).

- MATH (Hendrycks et al., 2021) & GSM8K (Cobbe et al., 2021): This dataset includes math questions with their correct answer included in the metadata. We can then verify the answer using the HF Math-Verify library (Kydlíček). We take the data as-is from `allenai/RLVR-MATH` and `allenai/RLVR-GSM`, and include them along with their Hebrew translations. The final score is either 0 (incorrect) or 1 (correct).

- Nikud (Diacriticization): This task is Hebrew-specific, where we provide the model with a sentence in Hebrew and expect the output with diacritics. This is an increasingly difficult task for LLMs to accomplish, as they both have to disambiguate the different meanings of each word, and also break down each word to the indivdiual letters as the diacritics are separate tokens. We use our in-house curated corpus of sentences diacritized by linguistic experts, providing us with a gold corpus to verify against. The final score is the percent of words correctly diacritized (between 0 and 1).

- Universal Dependencies (UD) Parsing: This task evaluates syntactic analysis by providing a sentence and expecting a full Universal Dependencies parse, including part-of-speech tags, dependency heads, and dependency relations. The model must correctly analyze the sentence and produce a correct by-word breakdown of the sentence. This is a very non-trivial task that even frontier-LLMs struggle

with. We use the data from the Hebrew Treebank (Zeldes et al., 2022; Tsarfaty, 2013; McDonald et al., 2013) as our gold standard (selecting only sentences which were not included in the earlier training phases). The final score is the percent of correct labels produced (between 0 and 1).

The final dataset consists of 46k examples.

### 5.3.2 Training

As with DPO training, we conducted the supervised fine-tuning on the same compute cluster on NVIDIA DGX Cloud Lepton, using the NVIDIA NeMo-RL Training Framework.

For this train we use the Megatron-LM backend for the training, and the vLLM backend (Kwon et al., 2023) for generation. We detail the full training parameters in Table 7.

We ran GRPO training on the `Dicta-LM 3.0` 24B and 1.7B thinking variants, with the model quickly learning. The reward graph for the 24B can be seen in Figure 4, and as can be seen the model quickly improves, and begins to plateau after 250 steps.

| Hyperparameter | Value |
|---|---|
| Micro batch size | 2 |
| Global batch size | 480 |
| Tensor Parallel (TP) | 4 |
| Max total sequence length | 8,192 |
| Num prompts per step | 64 |
| Num generations per prompt | 32 |
| KL penalty | 0.01 |
| KL type | k3 |
| Learning rate (LR) | 3e-7 |

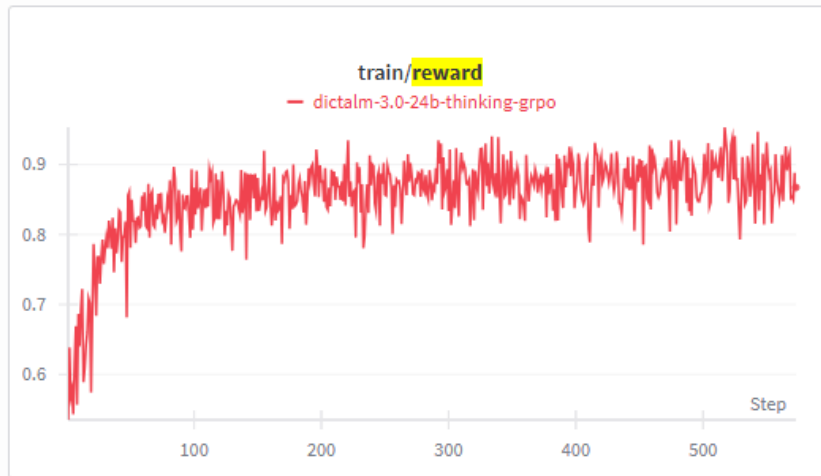Table 7: Training hyperparameters for the GRPO train.



Figure 4: Reward graph when running GRPO training on the 24B thinking variant.

## 6 Chat Model Evaluation

### 6.1 Hebrew Chat Evaluation

For Hebrew chat evaluation, we curated a set of five tasks to evaluate the performance of LLMs on Hebrew. The evaluation here differs from base model evaluation, where for the base model we provide few-shot examples, and here the model is provided with a 0-shot prompt. The tasks are also available in the form of a public leaderboard.[6] The five evaluation tasks are as follows:

---

[6]https://huggingface.co/spaces/hebrew-llm-leaderboard/chat-leaderboard

- **Summarization**: Evaluating the model's ability to summarize a document. The test corpus is comprised of 1,000 documents of various styles (news, wiki, etc.). We compare the test model's generated summary with a summary generated via Gemini-2.5-Pro, and check the win rate. If Gemini-2.5-Pro is rated higher, we assign a score of 0; if the test model's summary is better, we assign a score of 1. If the two are equally good, we assign a score of 0.5. We use GPT4o as the LLM as a Judge to determine the score here.

- **Translation**: The test corpus consists of a random set of English paragraphs, 20-40 words in length. For each one, we ask the test model to generate a translation to Hebrew, and we compare to the translation generated from Gemini-2.5-Pro, utilizing the same approach as detailed above with regard to summarization.

- **Israeli Trivia**: The test corpus consists of 300 trivia questions regarding Israel-related facts. We extend immense gratitude to Avraham Elitzur for helping us curate this dataset. We provide multiple-choice answers in the prompt and compute the accuracy score using exact string matching.

- **Winograd**: We use "A Translation of the Winograd Schema Challenge to Hebrew" (Shwartz, 2021). Models must resolve ambiguous pronouns using commonsense reasoning. We provide the ambiguous sentence and the question, with two choices for the model to choose from.

- **Nikud (Diacritization)**: We use a test corpus of sentences that were manually diacritized by expert linguists (ensuring that none were present in any phase of the training). Each sentence is scored according to the number of words were correctly diacritized. The final score is calculated as the macro average of the accuracies across the sentences in the corpus.

We present the results of our models, comparing them to other models of similar sizes in Table 8. As can be seen, the 24B model performs extremely well on language-specific tasks, which usually require a more subtle understanding of the language. In addition, the model also performs on par / better than models >30% larger than it on tasks which require more extensive world knowledge. The Nemotron-12B performs very well competitive with gemma-3-12b-it, yet with the major advantage that the Nemotron model is based on a hybrid-SSM architecture, allowing it to scale across longer context lengths with higher throughput and a lower memory footprint. The 1.7B model showcases remarkable performance for its size significantly outperforming other models in that size range.

| Model | Summarization | Translation | Winogrande | Trivia | Nikud |
|---|---|---|---|---|---|
| **DictaLM-3.0-24B-Thinking** | **56.86** | **30.09** | 78.06 | **60.13** | **86.86** |
| gemma3-27B-it | 44.54 | 26.73 | **79.86** | 45.51 | 60.21 |
| aya-expanse-32B | 29.46 | 17.10 | 80.58 | 53.82 | 45.40 |
| Llama-3.3-70B-Instruct | 37.83 | 19.31 | 83.45 | **60.13** | 4.05 |
| Smaller models (~12B) | | | | | |
| **DictaLM-3.0-Nemotron-12B-Instruct** | 33.27 | 13.50 | 73.74 | **45.18** | **76.12** |
| gemma-3-12b-it | **39.48** | **16.50** | **75.90** | 44.85 | 51.78 |
| Qwen3-14B (think) | 15.83 | 0.90 | 73.38 | 41.86 | 4.73 |
| Tiny models (~1.7B) | | | | | |
| DictaLM-3.0-1.7B-Instruct | 9.72 | 2.16 | **58.2** | 30.21 | **52.76** |
| **DictaLM-3.0-1.7B-Thinking** | **10.22** | **2.51** | 55.76 | **31.23** | 47.2 |
| gemma-3-1b-it | 0.35 | 0.15 | 47.84 | 26.58 | 3.44 |
| Qwen3-1.7B (think) | 0.4 | 0 | 51.08 | 21.59 | 2.93 |

Table 8: Evaluation scores for DictaLM-3.0 variants and baselines (higher is better).

## 6.2 English Chat Evaluation

For the English chat evaluation, we evaluate using the Olmes library (AllenAI, 2025) on the same suite of tasks as Olmo3 (Ai2, 2025)[7]. We present the results for the models in Tables 9 and 10. DictaLM 3.0 24B showcases very strong results across the board on the English benchmarks, specifically with regard to instruction following and chat, making it an ideal choice for downstream use. DictaLM 3.0 12B shows strong knowledge retention, and performing competitively on the mathematical benchmarks.

In Table 11 we compare the behavior of the instruct and thinking variants of the 24B model, analyzing the efficiency of the model versus the resulting accuracy gains. Similar to the findings of Shmidman et al. (2025), we observe that the thinking model delivers notable accuracy improvements without becoming overly verbose. In several cases, it even achieves a significant boost in accuracy while generating *fewer* tokens. We believe this occurs because the reasoning content is more compact and efficient than standard output, enabling the model to produce high-quality answers without incurring substantial inference costs.

| Category | Benchmark | DictaLM 3.0 24B-Think | Gemma 3 27B | Mistral Small 3.1 |
|---|---|---|---|---|
| Math | MATH | 86.41 | **87.40** | 67.93 |
| | OMEGA | **28.38** | 24.00 | 13.76 |
| Reasoning | BigBenchHard | 73.00 | **82.40** | 71.91 |
| | ZebraLogic | **47.80** | 24.80 | 23.10 |
| | AGI Eval (EN) | **82.93** | 76.90 | 75.87 |
| Inst. Following | IFEval | **88.17** | 85.40 | 79.48 |
| Knowledge | MMLU | **85.93** | 74.60 | 83.11 |
| | PopQA | **30.59** | 30.20 | 29.70 |
| | GPQA | **55.13** | 45.00 | 44.87 |
| Chat | AlpacaEval 2 LC | **74.11** | 65.50 | 41.22 |

Table 9: Performance comparison of mid-size models (24B–27B parameters). DictaLM-3.0-24B-Thinking achieves the strongest results on instruction following and chat, while Gemma 3 27B leads on several reasoning tasks.

| Category | Benchmark | Gemma 3 12B | DictaLM 3.0 12B-Inst |
|---|---|---|---|
| Math | MATH | **84.25** | 74.99 |
| | OMEGA | 15.19 | 15.19 |
| Reasoning | BigBenchHard | 70.10 | **71.38** |
| | ZebraLogic | 0.00 | **20.90** |
| | AGI Eval (EN) | **73.91** | 73.75 |
| IF | IFEval | **82.25** | 72.07 |
| Knowledge | MMLU | 77.91 | **80.16** |
| | PopQA | 22.64 | **26.31** |
| | GPQA | 32.58 | **38.61** |
| Chat | AlpacaEval 2 LC | **50.57** | 17.45 |

Table 10: Performance comparison of 12B-parameter models. DictaLM-3.0-12B-Instruct shows strong knowledge retention (MMLU, GPQA) while Gemma 3 12B excels at mathematical reasoning.

## 7 Conclusion

In this report we introduced Dicta-LM 3.0, a new suite of Hebrew-focused sovereign LLMs spanning 24B, 12B, and 1.7B parameters, each trained on large-scale Hebrew corpora and extended with long-context capabilities. Through continuous pre-training, supervised fine-tuning, and reinforcement-based post-training, the models achieve substantial gains across a wide range of Hebrew tasks, and maintain strong performance on English benchmarks despite the heavy Hebrew emphasis.

---

[7]We don't evaluate on the coding tasks, as we specifically did not focus on code generation abilities (see Section 5.1.1). In addition, we didn't evaluate on the following tasks since the datasets used in olmes aren't public as of the publication date: AIME 2024, AIME 2025, IFBench

| Category | Benchmark | Accuracy (%) | | | Tokens | | |
|---|---|---|---|---|---|---|---|
| | | Inst. | Think | Δ | Inst. | Think | Δ |
| Math | MATH | 60.83 | 86.41 | +25.6 | 2,233 | 1,457 | −35% |
| | OMEGA | 14.16 | 28.38 | +14.2 | 2,319 | 5,145 | +122% |
| Reasoning | BigBenchHard | 64.12 | 73.00 | +8.9 | 535 | 669 | +25% |
| | ZebraLogic | 23.60 | 47.80 | +24.2 | 6,090 | 5,495 | −10% |
| | AGI Eval (EN) | 73.12 | 82.93 | +9.8 | 650 | 994 | +53% |
| Inst. Following | IFEval | 87.80 | 88.17 | +0.4 | 666 | 2,230 | +235% |
| Knowledge | MMLU | 81.11 | 85.93 | +4.8 | 294 | 398 | +35% |
| | PopQA | 27.83 | 30.59 | +2.8 | 33 | 794 | +2306% |
| | GPQA | 41.96 | 55.13 | +13.2 | 1,481 | 2,331 | +57% |
| Chat | AlpacaEval 2 LC | 31.11 | 74.11 | +43.0 | 475 | 1,630 | +243% |
| **Average** | | – | – | +15.7 | 1,478 | 2,114 | +303% |

Table 11: Efficiency analysis: DictaLM-3.0-24B Instruct vs Thinking variants. The Thinking model achieves substantial accuracy gains across most benchmarks. Negative token overhead (green) indicates the Thinking model used *fewer* tokens while achieving better accuracy.

To address gaps in evaluating non-English chat models, we proposed a new Hebrew chat benchmark suite covering summarization, translation, reasoning, Israeli knowledge, and diacritization. The resulting evaluations show clear improvements over the base models and competitive performance against significantly larger multilingual LLMs. The 24B model, in particular, sets a new standard for Hebrew-capable open-weight models and demonstrates that sovereign LLMs can outperform generalist models when trained with focused data and targeted post-training.

Our work provides concrete methodology - data construction, training strategy, and evaluation tools - for adapting frontier LLMs to low-resource languages. We hope this effort supports broader development of sovereign language models and encourages further progress in multilingual NLP.

We are happy to release the models to the public with a permissive license:

- DictaLM-3.0-24B-Base - 24B Base model

- DictaLM-3.0-24B-Thinking - 24B chat model with reasoning and tool-calling capabilities.

- DictaLM-3.0-Nemotron-12B-Base - 12B base model.

- DictaLM-3.0-Nemotron-12B-Instruct - 12B chat model with tool-calling capabilities.

- DictaLM-3.0-1.7B-Base - 1.7B Base model

- DictaLM-3.0-1.7B-Instruct - 1.7B chat model with tool-calling capabilities.

- DictaLM-3.0-1.7B-Thinking - 1.7B chat model with reasoning and tool-calling capabilities.

## Acknowledgments

# References

Ai2. 2025. Olmo 3: Charting a path through the model flow to lead open-source ai. Blog post. Allen Institute for AI (AI2).

AI@Meta. 2024. Llama 3 model card.

AllenAI. 2025. OLMES: Open language model evaluation system. https://github.com/allenai/olmes. GitHub repository.

Janisha Anand and Sowmyan Soman. 2025. Introducing nvidia dgx cloud lepton: A unified ai platform built for developers. Blog post, NVIDIA Developer Blog.

Anthropic. 2024. The claude 3 model family: Opus, sonnet, haiku. *Papers with Code*.

Elie Bakouch, Loubna Ben Allal, Anton Lozhkov, Nouamane Tazi, Lewis Tunstall, Carlos Miguel Patiño, Edward Beeching, Aymeric Roucher, Aksel Joonas Reedi, Quentin Gallouédec, Kashif Rasul, Nathan Habib, Clémentine Fourrier, Hynek Kydlicek, Guilherme Penedo, Hugo Larcher, Mathieu Morlon, Vaibhav Srivastav, Joshua Lochner, Xuan-Son Nguyen, Colin Raffel, Leandro von Werra, and Thomas Wolf. 2025. SmolLM3: smol, multilingual, long-context reasoner. https://huggingface.co/blog/smollm3.

Ben-Yehuda. 2024. Ben yehuda project.

Akhiad Bercovich, Itay Levy, Izik Golan, Mohammad Dabbah, Ran El-Yaniv, Omri Puny, Ido Galil, Zach Moshe, Tomer Ronen, Najeeb Nabwani, Ido Shahaf, Oren Tropp, Ehud Karpas, Ran Zilberstein, Jiaqi Zeng, Soumye Singhal, Alexander Bukharin, Yian Zhang, Tugrul Konuk, Gerald Shen, Ameya Sunil Mahabaleshwarkar, Bilal Kartal, Yoshi Suhara, Olivier Delalleau, Zijia Chen, Zhilin Wang, David Mosallanezhad, Adi Renduchintala, Haifeng Qian, Dima Rekesh, Fei Jia, Somshubra Majumdar, and Vahid Noroozi et al. 2025. Llama-nemotron: Efficient reasoning models.

Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners.

Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. 2018. Think you have solved question answering? try arc, the ai2 reasoning challenge.

Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*.

DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, Aixin Liu, Bing Xue, Bingxuan Wang, Bochao Wu, Bei Feng, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng Wang, Honghui Ding, Huajian Xin, Huazuo Gao, Hui Qu, Hui Li, Jianzhong Guo, Jiashi Li, Jiawei Wang, Jingchang Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, J. L. Cai, Jiaqi Ni, Jian Liang, Jin Chen, and Kai Dong et al. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning.

Thomas Mesnard Gemma Team, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, Pouya Tafti, Léonard Hussenot, and et al. 2024. Gemma.

Dirk Groeneveld, Iz Beltagy, Pete Walsh, Akshita Bhagia, Rodney Kinney, Oyvind Tafjord, Ananya Harsh Jha, Hamish Ivison, Ian Magnusson, Yizhong Wang, Shane Arora, David Atkinson, Russell Arthur, Khyathi Raghavi Chandu, Arman Cohan, Jennifer Dumas, Yanai Elazar, Yuling Gu, Jack Hessel, Tushar Khot, William Merrill, Jacob Morrison, Niklas Muennighoff, Aakanksha Naik, Crystal Nam, Matthew E. Peters, Valentina Pyatkin, Abhilasha Ravichander, Dustin Schwenk, Saurabh Shah, Will Smith, Emma Strubell, Nishant Subramani, Mitchell Wortsman, Pradeep Dasigi, Nathan Lambert, Kyle Richardson, Luke Zettlemoyer, Jesse Dodge, Kyle Lo, Luca Soldaini, Noah A. Smith, and Hannaneh Hajishirzi. 2024. Olmo: Accelerating the science of language models. *arXiv preprint arXiv:2402.00838*.

Nathan Habib, Clémentine Fourrier, Hynek Kydlíček, Thomas Wolf, and Lewis Tunstall. 2023. Lighteval: A lightweight framework for llm evaluation.

Eric Harper, Somshubra Majumdar, Oleksii Kuchaiev, Li Jason, Yang Zhang, Evelina Bakhturina, Vahid Noroozi, Sandeep Subramanian, Koluguri Nithin, Huang Jocelyn, Fei Jia, Jagadeesh Balam, Xuesong Yang, Micha Livne, Yi Dong, Sean Naren, and Boris Ginsburg. NeMo: a toolkit for Conversational AI and Large Language Models.

hebrew-llm-leaderboard. 2025. Hebrew llm leaderboard. `https://huggingface.co/spaces/hebrew-llm-leaderboard/leaderboard`. Hugging Face Spaces.

Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. 2021. Measuring mathematical problem solving with the math dataset. *NeurIPS*.

Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. 2023. Efficient memory management for large language model serving with pagedattention. In *Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles*.

Hynek Kydlíček. Math-Verify: Math Verification Library.

Nathan Lambert, Jacob Morrison, Valentina Pyatkin, Shengyi Huang, Hamish Ivison, Faeze Brahman, Lester James V. Miranda, Alisa Liu, Nouha Dziri, Shane Lyu, Yuling Gu, Saumya Malik, Victoria Graf, Jena D. Hwang, Jiangjiang Yang, Ronan Le Bras, Oyvind Tafjord, Chris Wilhelm, Luca Soldaini, Noah A. Smith, Yizhong Wang, Pradeep Dasigi, and Hannaneh Hajishirzi. 2024. Tulu 3: Pushing frontiers in open language model post-training.

Yifei Liu, Li Lyna Zhang, Yi Zhu, Bingcheng Dong, Xudong Zhou, Ning Shang, Fan Yang, and Mao Yang. 2025. rstar-coder: Scaling competitive code reasoning with a large-scale verified dataset.

Anton Lozhkov, Loubna Ben Allal, Leandro von Werra, and Thomas Wolf. 2024. Fineweb-edu: the finest collection of educational content.

Ryan T McDonald, Joakim Nivre, Yvonne Quirmbach-Brundage, Yoav Goldberg, Dipanjan Das, Kuzman Ganchev, Keith B Hall, Slav Petrov, Hao Zhang, Oscar Täckström, et al. 2013. Universal dependency annotation for multilingual parsing. In *Proc. of ACL*.

Mistral AI. 2025. Mistral small 3.1. `https://mistral.ai/news/mistral-small-3-1`. 24 B-parameter open-source multimodal model, Apache 2.0 license.

Dhruv Nathawani, Igor Gitman, Somshubra Majumdar, Evelina Bakhturina, Ameya Sunil Mahabaleshwarkar, , Jian Zhang, and Jane Polak Scowcroft. 2025. Nemotron-Post-Training-Dataset-v1.

NousResearch. 2024. Hermes-3 dataset. Hugging Face Dataset.

NVIDIA, :, Aarti Basant, Abhijit Khairnar, Abhijit Paithankar, Abhinav Khattar, Adithya Renduchintala, Aditya Malte, Akhiad Bercovich, Akshay Hazare, Alejandra Rico, Aleksander Ficek, Alex Kondratenko, Alex Shaposhnikov, Alexander Bukharin, Ali Taghibakhshi, Amelia Barton, Ameya Sunil Mahabaleshwarkar, Amy Shen, Andrew Tao, Ann Guan, Anna Shors, Anubhav Mandarwal, Arham Mehta, Arun Venkatesan, Ashton Sharabiani, Ashwath Aithal, Ashwin Poojary, and Ayush Dattagupta et al. 2025. Nvidia nemotron nano 2: An accurate and efficient hybrid mamba-transformer reasoning model.

Nvidia. 2025. NeMo RL: A scalable and efficient post-training library. `https://github.com/NVIDIA-NeMo/RL`. GitHub repository.

OpenAI, :, Aaron Hurst, Adam Lerer, Adam P. Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, Aleksander Mądry, Alex Baker-Whitcomb, Alex Beutel, Alex Borzunov, Alex Carney, Alex Chow, Alex Kirillov, Alex Nichol, Alex Paino, Alex Renzin, Alex Tachard Passos, Alexander Kirillov, Alexi Christakis, and Alexis Conneau et al. 2024a. Gpt-4o system card.

OpenAI. 2025. gpt-oss-120b & gpt-oss-20b model card.

OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, et al.Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian, Jeff Belgum, and et al. 2024b. Gpt-4 technical report.

Pedro Javier Ortiz Suárez, Benoît Sagot, and Laurent Romary. 2019. Asynchronous pipelines for processing huge corpora on medium to low resource infrastructures. Proceedings of the Workshop on Challenges in the Management of Large Corpora (CMLC-7) 2019. Cardiff, 22nd July 2019, pages 9 – 16, Mannheim. Leibniz-Institut für Deutsche Sprache.

Guilherme Penedo, Hynek Kydlíček, Vinko Sabolčec, Bettina Messmer, Negar Foroutan, Amir Hossein Kargaran, Colin Raffel, Martin Jaggi, Leandro Von Werra, and Thomas Wolf. 2025. Fineweb2: One pipeline to scale them all – adapting pre-training data processing to every language.

Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D. Manning, and Chelsea Finn. 2023. Direct preference optimization: Your language model is secretly a reward model.

Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2019. Exploring the limits of transfer learning with a unified text-to-text transformer. *arXiv e-prints*.

Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. 2019. Winogrande: An adversarial winograd schema challenge at scale.

Holger Schwenk, Guillaume Wenzek, Sergey Edunov, Edouard Grave, and Armand Joulin. 2020. Ccmatrix: Mining billions of high-quality parallel sentences on the web.

Sefaria. 2024. Sefaria: A living library of jewish texts.

Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, Y. K. Li, Y. Wu, and Daya Guo. 2024. Deepseekmath: Pushing the limits of mathematical reasoning in open language models.

Shaltiel Shmidman, Asher Fredman, Oleg Sudakov, and Meriem Bendris. 2025. Learning to reason: Training llms with gpt-oss or deepseek r1 reasoning traces.

Shaltiel Shmidman, Avi Shmidman, Amir DN Cohen, and Moshe Koppel. 2024. Adapting llms to hebrew: Unveiling dictalm 2.0 with enhanced vocabulary and instruction capabilities.

Mohammad Shoeybi, Mostofa Patwary, Raul Puri, Patrick LeGresley, Jared Casper, and Bryan Catanzaro. 2020. Megatron-LM: Training multi-billion parameter language models using model parallelism.

Vered Shwartz. 2021. A translation of the winograd schema challenge to hebrew.

Daria Soboleva, Faisal Al-Khateeb, Robert Myers, Jacob R Steeves, Joel Hestness, and Nolan Dey. 2023. SlimPajama: A 627B token cleaned and deduplicated version of RedPajama. https://www.cerebras.net/blog/slimpajama-a-627b-token-cleaned-and-deduplicated-version-of-redpajama.

Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. 2019. CommonsenseQA: A question answering challenge targeting commonsense knowledge. 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 4149–4158, Minneapolis, Minnesota. Association for Computational Linguistics.

Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M. Dai, Anja Hauth, and et al. 2024. Gemini: A family of highly capable multimodal models.

Ryan Teknium, Roger Jin, Jai Suphavadeeprasit, Dakota Mahan, Jeffrey Quesnelle, Joe Li, Chen Guang, Shannon Sands, and Karan Malhotra. 2024. Hermes 3 technical report. Function-calling schema described — tool definitions in <tools>, calls in <tool_call>/<tool_response>.

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan

Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. Llama 2: Open foundation and fine-tuned chat models.

Reut Tsarfaty. 2013. A unified morpho-syntactic scheme of stanford dependencies. In *Proc. of ACL*.

Zhangchen Xu, Adriana Meza Soria, Shawn Tan, Anurag Roy, Ashish Sunil Agrawal, Radha Poovendran, and Rameswar Panda. 2025. Toucan: Synthesizing 1.5m tool-agentic data from real-world mcp environments.

An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, Chujie Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng Hu, Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jing Zhou, Jingren Zhou, Junyang Lin, Kai Dang, Keqin Bao, Kexin Yang, Le Yu, Lianghao Deng, Mei Li, Mingfeng Xue, Mingze Li, Pei Zhang, Peng Wang, Qin Zhu, Rui Men, Ruize Gao, Shixuan Liu, Shuang Luo, Tianhao Li, Tianyi Tang, Wenbiao Yin, Xingzhang Ren, Xinyu Wang, Xinyu Zhang, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yinger Zhang, Yu Wan, Yuqiong Liu, Zekun Wang, Zeyu Cui, Zhenru Zhang, Zhipeng Zhou, and Zihan Qiu. 2025. Qwen3 technical report.

Amir Zeldes, Nick Howell, Noam Ordan, and Yifat Ben Moshe. 2022. A second wave of UD Hebrew treebanking and cross-domain parsing. In *Proceedings of EMNLP 2022*, pages 4331–4344, Abu Dhabi, UAE.

Jeffrey Zhou, Tianjian Lu, Swaroop Mishra, Siddhartha Brahma, Sujoy Basu, Yi Luan, Denny Zhou, and Le Hou. 2023. Instruction-following evaluation for large language models.

## A   Appendix: Lighteval command for base model evaluation on English tasks

```
lighteval vllm "model_name=google/gemma-3-27b-pt,tensor_parallel_size=2,max_model_length=32768,
    trust_remote_code=True" "helm|commonsenseqa|5,leaderboard|arc:challenge|25,leaderboard|
    winogrande|5"
```