

Cotype Nano

Agenda

01 LLMs & SMLs

02 1 stage SFT

03 2 stage SFT

04 Datasets

05 Evaluation

06 Quantization

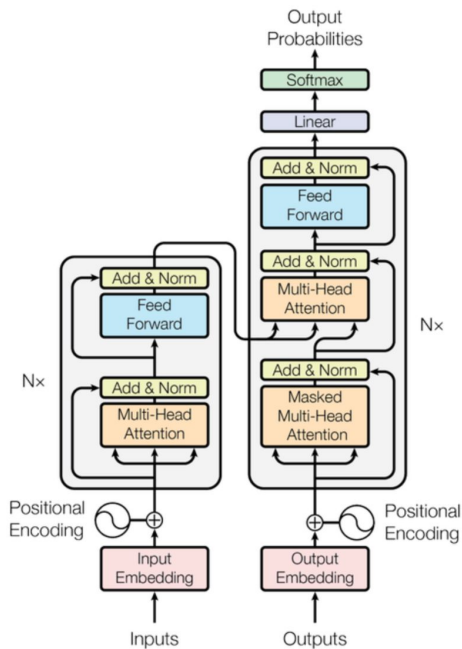
01

LLMs and SLMs

LLMs and SLMs

BERT

Encoder



GPT

Decoder

Output
Probabilities

Softmax

Linear

Add & Norm

Feed
Forward

Add & Norm

Multi-Head
Attention

Add & Norm

Masked
Multi-Head
Attention

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

Add & Norm

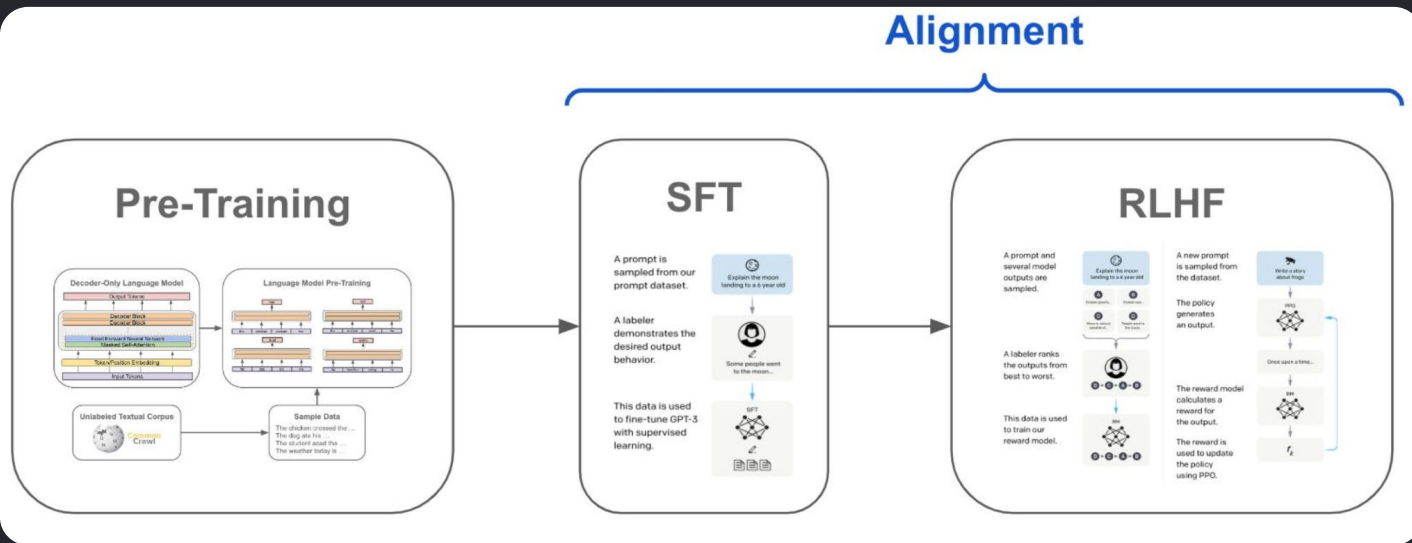
Add & Norm

Add & Norm

02

One-stage SFT 🤗

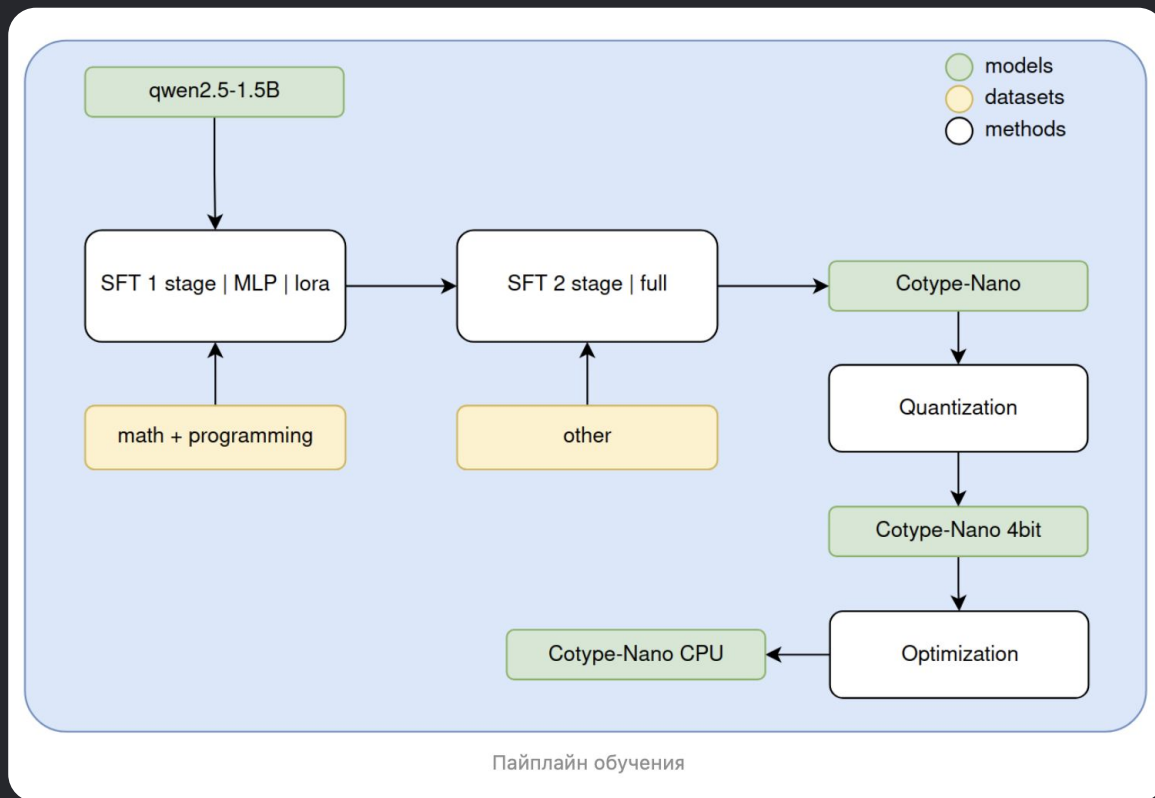
Alignment



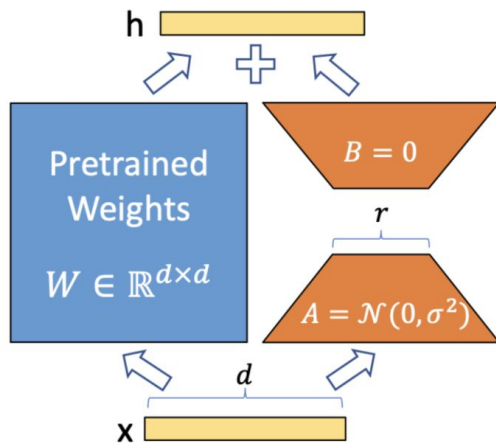
03

Two-stage SFT 🤗

Two-stage SFT



Two-stage SFT. Stage 1



Datasets:

Math (200k)

Programming (~50k)

Two-stage SFT. Stage 2

Full model

Datasets:

General instruction
datasets

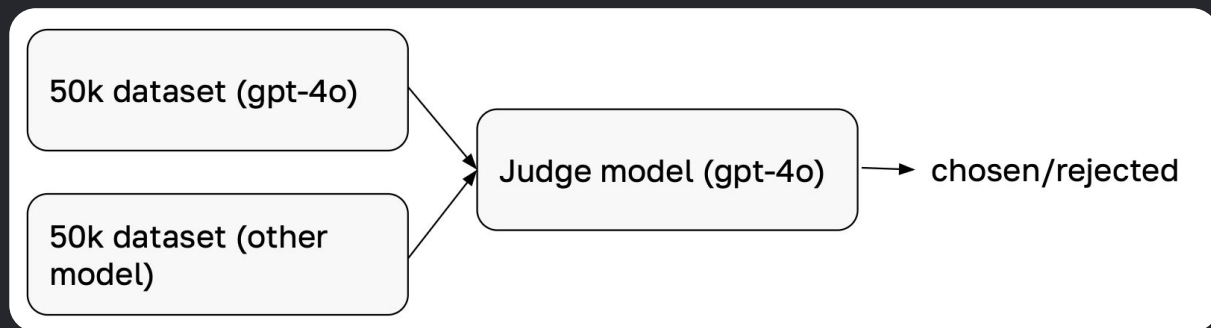
04

Dataset

 [microsoft](#)/**orca-math-word-problems-200k**

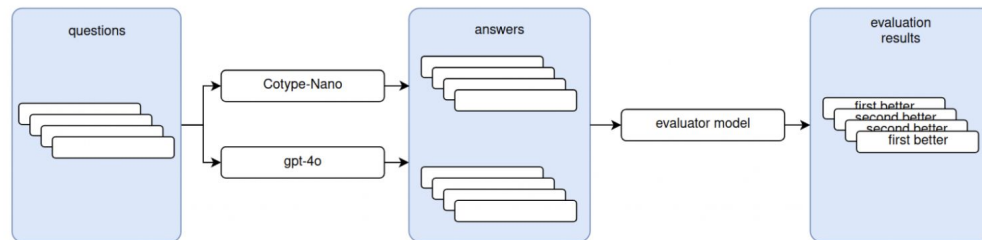
Internal programming dataset

Internal INST dataset + part of open source datasets

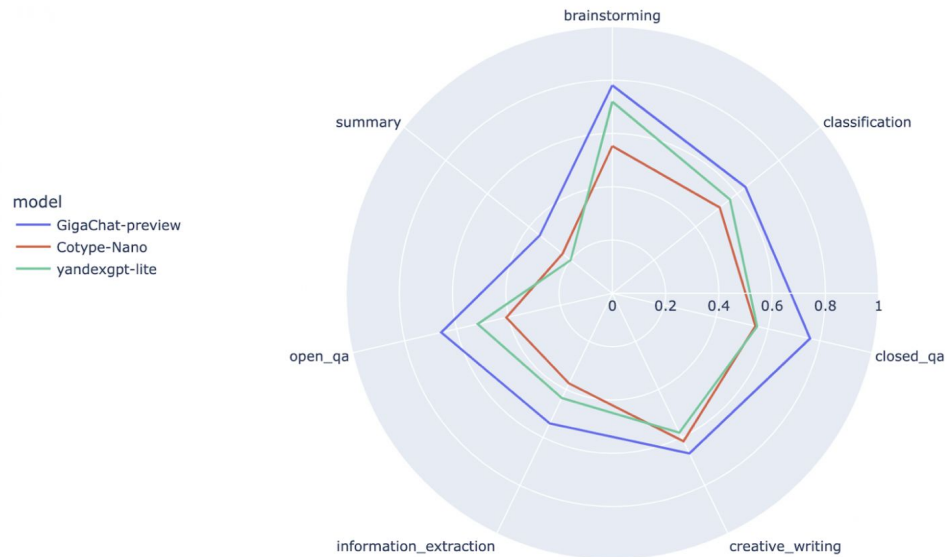


05

Evaluation 



autoSBS оценка для моделей



Локальная автооценка с помощью gpt-4

Model	Score	95% CI	Avg. #Tokens
Cotype-Nano	30.2	(-1.3, 2.2)	542
vikhr-it-5.3-fp16-32k	27.8	(-2.1, 1.5)	519
vikhr-it-5.3-fp16	22.7	(-1.7, 1.8)	523
Cotype-Nano 4bit	22.5	(-1.4, 2.1)	582
storm-7b	20.6	(-1.6, 1.7)	419
Vikhrmodels-Vikhr-Llama-3.2-1B-instruct	19.4	(-1.6, 1.3)	958
neural-chat-7b-v3-3	19.0	(-1.6, 1.6)	927
gigachat_lite	17.2	(-1.4, 1.7)	276
Vikhrmodels-vikhr-qwen-1.5b-it	13.1	(-1.1, 1.1)	2495
meta-llama-Llama-3.2-1B-Instruct	4.0	(-0.6, 0.8)	1240

Ru llm arena

MTSAIR-Cotype-Nano-Uncensored	50.51
sfr-iterative-dpo-llama-3-8b-r	50.06
gpt-3.5-turbo-0125	50
glm-4-9b-chat	49.75
c4ai-command-r-v01	48.95
-ruadapt_llama3_extended_gm_ft_v5d1	48.16
-kolibri-mistral-0427-upd	47.91
-ruadapt_llama3_8b_instruct_extended_lep_ft-externa	47.81
MTSAIR-Cotype-Nano	47.74
llama-3-instruct-8b-sppo-iter3	47.45
-ruadapt_saiga_v7_lep_ft_external_infer	47.36
Vikhrmodels-Vikhr-Qwen-2.5-1.5b-Instruct	47.23

Ru llm arena

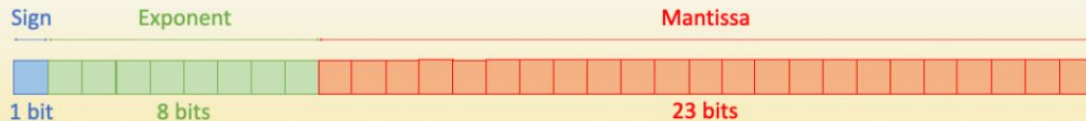
MTSAIR-Cotype-Nano-1B-v2	36.77
aya-23-8b	36.26
meta-llama-3-8b-instruct	35.07
openchat-3.5-0106	33.79
mistral-7b-instruct-v0.3	32.92
vikhr-it-5.2-fp16-cp	31.74
gigachat_pro	31.37
Vikhrmodels-Vikhr-Llama-3.2-1B-instruct	19.04
gigachat_lite	17.2
Vikhrmodels-Vikhr-Qwen-2.5-0.5b-Instruct	16.5

06

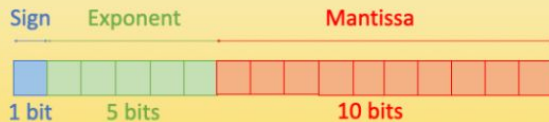
Quantization

Pre-training quantization

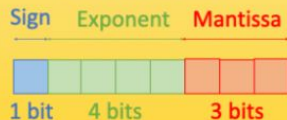
Float 32



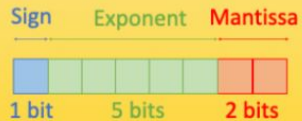
Float 16



Float 8 E4M3



Float 8 E5M2



Post-training quantization



TinyChat Computer
(Jetson Orin Nano)



Raspberry Pi
(ARM CPU)

fp16



int4



AWQ



Quantization Algorithm: AWQ

Inference System: TinyChat



MacBook
(Apple M1)



AI PC
(CPU / GPU)

Post-training quantization

 MTSAIR/**MultiVerse_70B_AWQ** 

 MTSAIR/**multi_verse_model_GPTQ** 

Thank you for your attention