# Llama 3.1 Paper

Pretraining :- ① filtering of large scale training corpus

② Devlopment of model Architecture
   ↳ Scaling laws for model size

③ efficient large scale pretraining
   technique development

④ pretraining recipe development


Dataset : end of 2023

Removed domains that has personal info of People
and adult content, insafe contents

Used HTML for high quality diverse tent.

built custom parser, that extracts HTML
boiler plate removal, content recd

Used HTML ⟶ tent extractor like.
   [jina.ai]

found ⟶ Markdown is harmful to
model performance. that
is trained on web date
than plane tent.


Several rounds of deduplication

Heuristic filtering : → to remove low quit quality
document

Multilingual Data :→ followed Similar pipeline
Added filter to remove personal data

→ language identification using fast text

→ document & line level deduplication

→ Applied language level heuristics } remove
model based filtering } low quality
document

Data Mix

Knowledge classification → Categorise types of info
↳ down sample data Categories
that are over represented

Scaling laws :-

## llama 3.1 paper (continuation)

→ to reduce memory cost → deallocate the Tensors that will not be used for further computation, including input & output tensors for each pipeline

[How to figure out which tensors are not being used?]

→
[Pg 12] Developed memory consumption estimator and performance projection tool — to explore parallelism config and overall training performance and identify memory gap

Training :-
① initial pre training —— Adam, $\alpha = 8 \times 10^{-5}$
② long content pretraining — 800 B tokens, 128 K content
③ Annealing

# llama 3.1 paper (continuation)

→ to reduce memory cost → deallocate the Tensors that will not be used for further computation including input & output tensors in each pipeline

→ Developed memory consumption estimator and performance projection tool — to explore parallelism config and overall training performance and identify memory gap

[Pg 12]

Training :-
① initial pre training ——— Adam, $\alpha = 8 \times 10^{-5}$
② long content pretraining — 800 B tokens, 128 k content
③ Annealing