# Hallucinations in LLM

generated content ⇒ Non sensical or
unfaithful to the provided source content

Intrinsic
Hallucination

Extrinsic
Hallucination

## Survey :

factuality
Hallucination

generated content &
verifiable
Real world facts
typically maniferting
a factual inconsistency
or fabrication

faithfulness
Hallucination

divergence of
generated content
from user's
information.
consistency

Instruction

Logical

factual
inconsistency

factual
fabrication

Context

## Understand what LLM is all about :

LLM

Series of transformer Based
Model

Large corpora

By scaling the
size of the
Model & the
Capability

of training data

LLM

In Context          Chain        Instruction
Learning           Thought       following
[ICL]              Reasoning

Training Stages of LLM:

① Pre training

② Supervised fine tuning

③ RLHFC (Reinforcement learning with Human feedback)

① PRETRAINING: — MORE IMPORTANT
                [acquire knowledge & skills.]
Through Step1: LLM aims to predict
the next word auto regressively
knowledge:

    ① Language Syntax
    ② World knowledge
    ③ Reasoning abilities and
  providing robust found ahone for
subsequent fine tuning tasks.

Essence of (PRETRAINING) ⟹ predicts the
probability of the next word in the
sequence.

② SUPERVISED FINE TUNING:

[PRETRAINING] primarily optimizes for
                                   Completion.

PRETRAINED LLM ⟹ Completion Machine.

LEADS to (MISALIGN) between the next word prediction objective of the LLM and the user's objective of obtaining Desired Responses

(TO BRIDGE this GAP) SFT Introduced.

which involves further training LLM's using a meticulously annotated set of [INSTRUCTION, RESPONSE] pairs.

③ RLHF

human preferences ⇌ LLM

More are lots of Methods But (RLHF) stands out as an Institute solution.

RHF → [PREFERENCE MODEL] trained to predict preference ranking given [PROMPT + PAIR OF HUMAN LABELED RESPONSES]

(RLHF) optimize (LLM) → to generate output

(PREFERENCE MODEL) maximize the Reward provided by the typically employ (PPO) Proximal policy optimization