# MIRASOL 3B
## Multimodal Auto Regressive Model.

Main Challenges = Combining Multiple Heterogeneous Modalities
(Video) (Audio) and (Text)

[Video |Audio] ⟹ higher rate than Text
are roughly aligned in time.

(Video| Audio)
They are not synchronized with Text which comes as a global context.
(eg:) Title/ Des caption.

[Video| Audio Inputs are much larger Volumes] ∝ Video length increases

⟱

Needs more Compute time for these Modalities.

(Mirasol 3B) ⟹ ① Auto regressive Component
The syn chronized Modalities
( Audio / Video )
Component          non time
which are not       ⟩ aligned
with time but
‖ Sequential ‖

⟱

② Autoregressive
Context modalities
[ necessarily    aligned

To address the long sequence of Video/ Audio inputs ⟹ partition of video/ Audio
Sequences in [ Consecutive Snippets ] and
autoregressively process their representation.
‖ Combiner Mechanism ‖        audio /video

information jointly producing compact but expressive representations.

This makes the Model to take 512 frame input without increase in the Model Parameters.

[Note: Architecture for Video Language understanding commonly use a joint transformer [Video input + Text tokens] are processed auto regressively]

Inputs:- Input video sequence of $N$ frames

$$V = (v_1 f, v_2 f, \dots v_N f)$$

Audio waves signal of $M$ timesteps

$$a = (a_1 f, a_2 f, \dots a_M f)$$

$$t = \{t_1 f, t_2 f, \dots t_P f\} \quad \boxed{\text{Text Sequence}}$$

Partioning the video & Audio inputs:

$$\left\{ v_1 f, v_2 f, \dots v_k f \right\}, \left\{ v_{k+1} f, v_{k+2} f, \dots v_{2k} f \right\}$$

$$\underbrace{\qquad}_{V_1} \qquad \underbrace{\qquad}_{V_2} \dots$$

chunks

$$\left( v^f_{(T-1)k+1}, \dots v_N f \right)$$

$$\underbrace{\qquad}_{V_T}$$

$$\boxed{k = N \mid T}$$

$T = $ Non Overallipity chunks.

[Each of every chunk] ⟹ Latent features

L $\left[\begin{array}{l}\text{Video} - \text{spatio-data representation.}\\ \text{Extract sparse 3D tubes.}\end{array}\right]$ with

standard 2D patches. $\left.\right]$ are processed

using the ViT encoder

Audio — Represented as $\boxed{\text{spectrograms}}$

## Combiner Module:-

1) Combine video / Audio features at specific snippet at time ( ) joint representato

2) effe compress the representation from each [audio / video] snippet, allows the model to scale to longer videos.

$\hat{V} = \{ \hat{v_1}, \hat{v_2} \ldots \hat{v_N} \} \quad \Rightarrow$ video

$\hat{a} = \{ \hat{a_1}, \hat{a_2}, \ldots \hat{a_N} \} \Rightarrow$ audio

Composed of f features of size d (f, d) shape

Composed of s features fixed shape $\{s, d\}$

Combiner ① $U = \{ u_1, u_2, v_3, \ldots u_T \}$

Where $\boxed{u_t : \{ \hat{v}_t, \hat{a}_t \}}$ and

size $\{n, d\}$ where $n = f + s$

② Then maps $u$ to Lower dimensional features pace $x = \{ x_1, x_2, \ldots x_T \}$

Where $x_t$ has the shape $\{m, d\}$

Where $n >> m$

Combothers Two different
Architectures
① Standard Transformer the
② Token Tuning Machine.

Time Aligned Video Audio Auto regressive Moo

Condition ( Audio / Video ) representatibne.
form previous time intervals

$\mathcal{Z}_t$ ⇒ passed sequentially to auto
regrective Model.
T

$$p(v, a) = \prod_{i=1} \left[ p(v_{t+1}, a_{t+1} | h_t) \right]$$
Modality
reconstruction Model.

$$\left[ p(h_t | \mathcal{Z}_t) \right] \left[ p(\mathcal{Z}_t | v_t, a_t) \right]$$
Combiner.

⇓

Seshuated by latent Causal Model.