

# Emotional FusionBrain 4.0



26.09.2024

Irina Abdullaeva



Multimodality +  
Socialization =  
AI Assistant

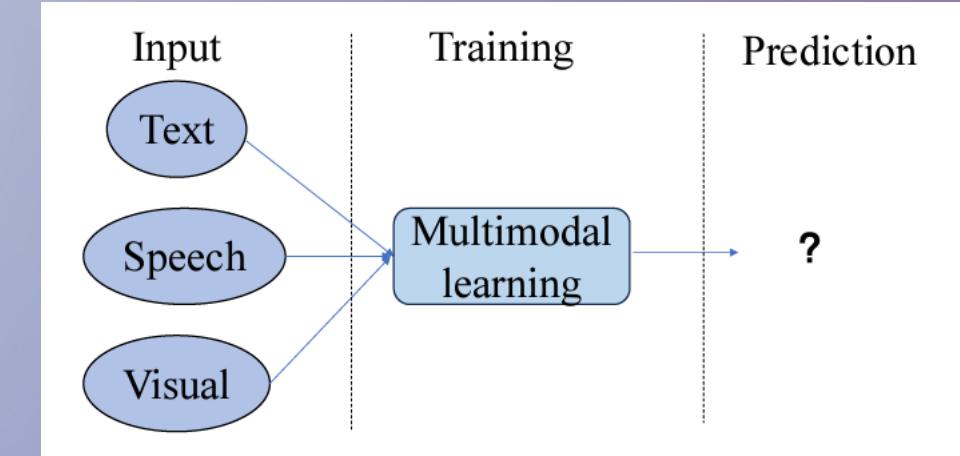
# Why we need multimodality?

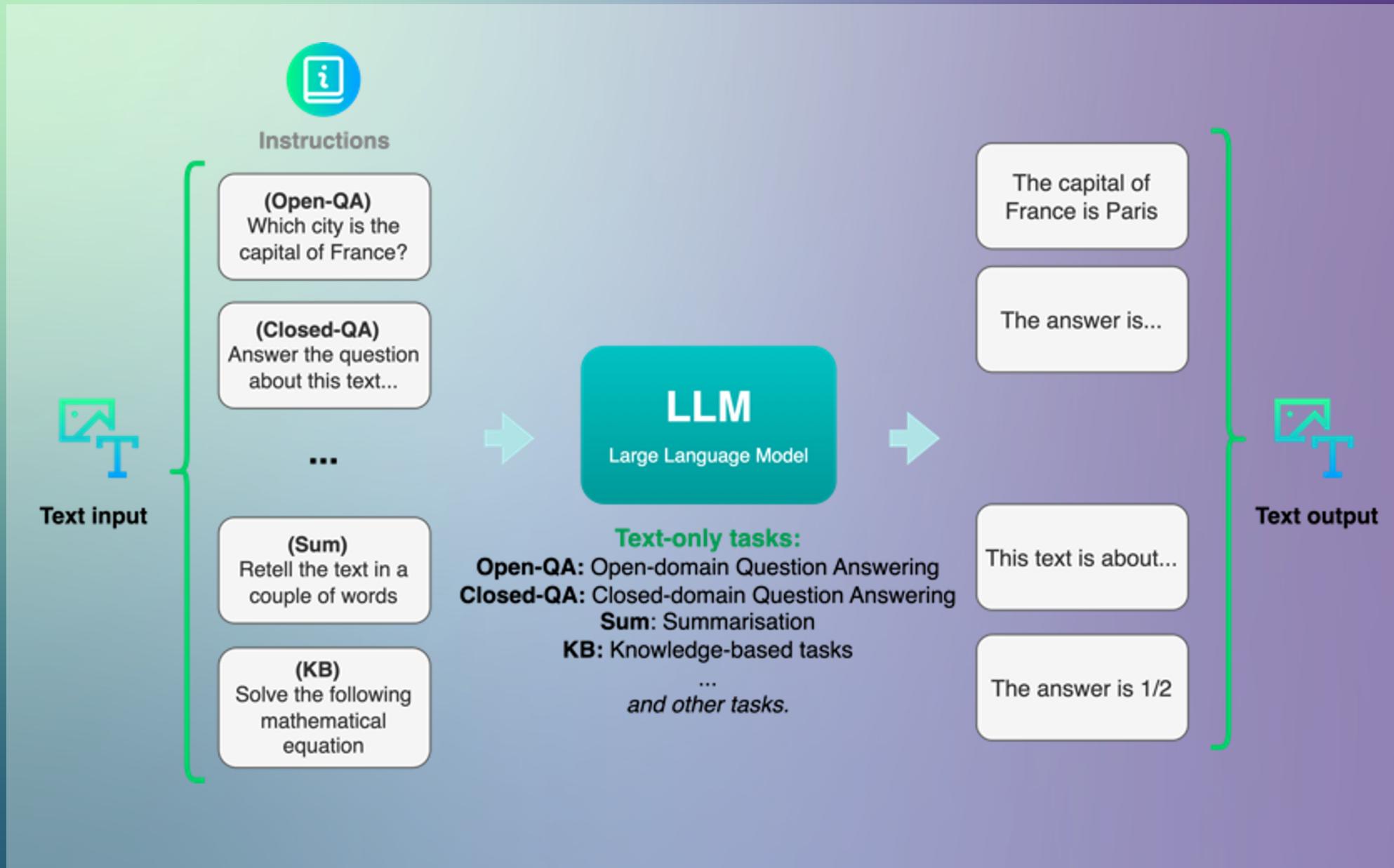
## Why is LLM popular?

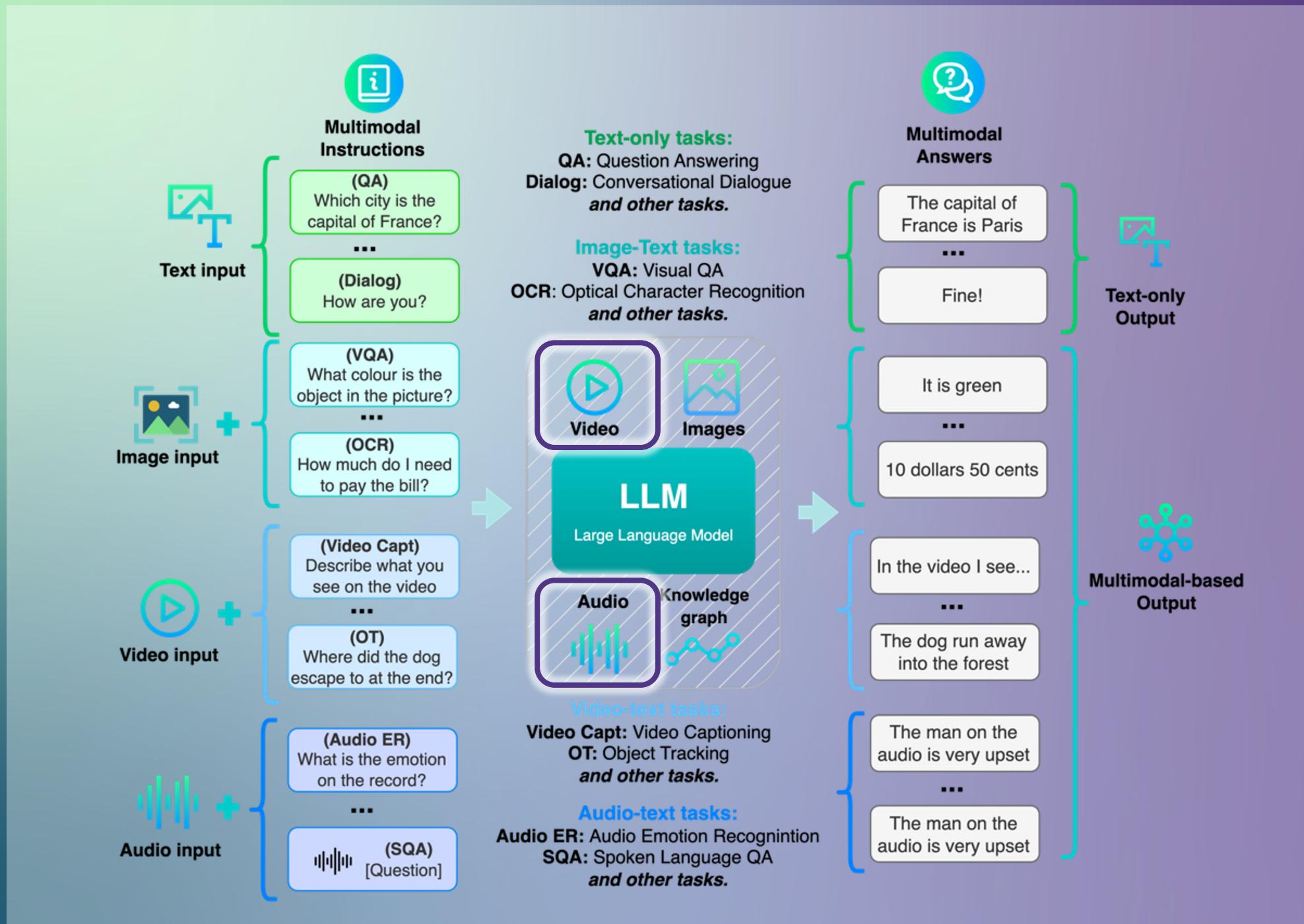
- Text – basis of communication
- Text – main source of the world knowledge
- Chatbots primarily need text

## Why do we need something more than text?

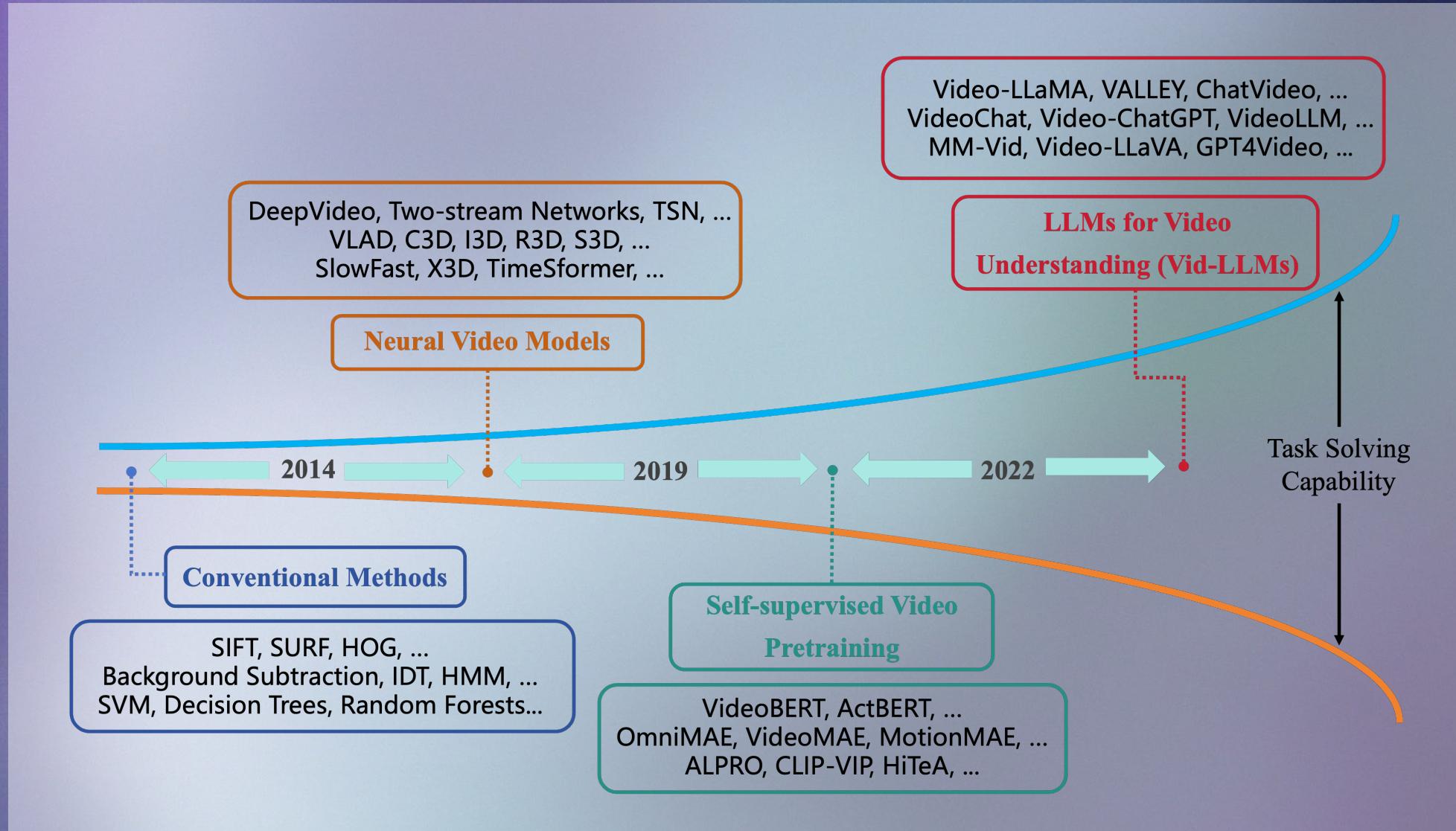
- People are used to communicate using **more than just text**
- To design a high-quality AI assistant, you need to extend perception abilities far beyond text
- Text + images + sound + **video** + ... — this is the perfect NATURAL scenario for interaction



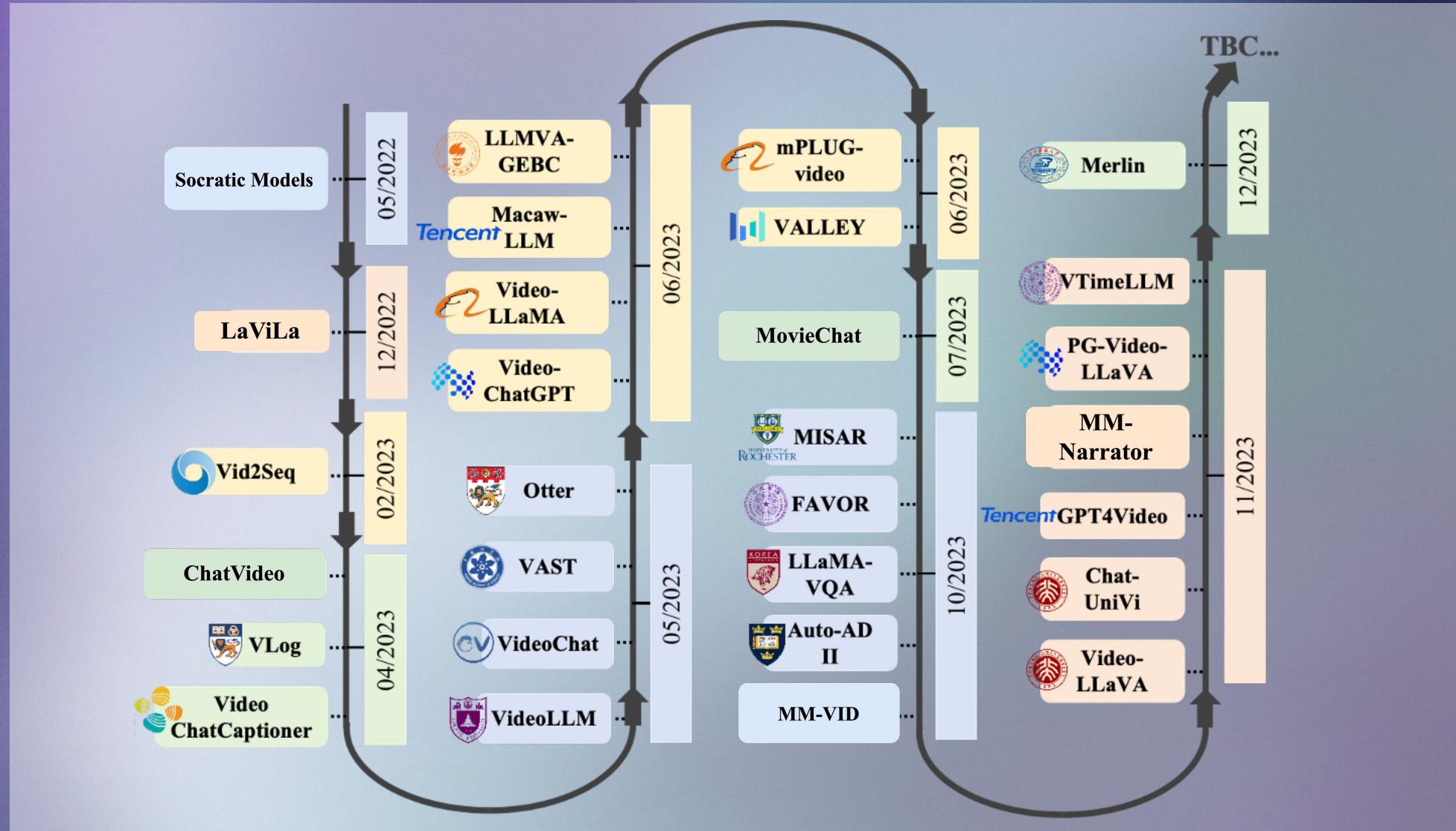




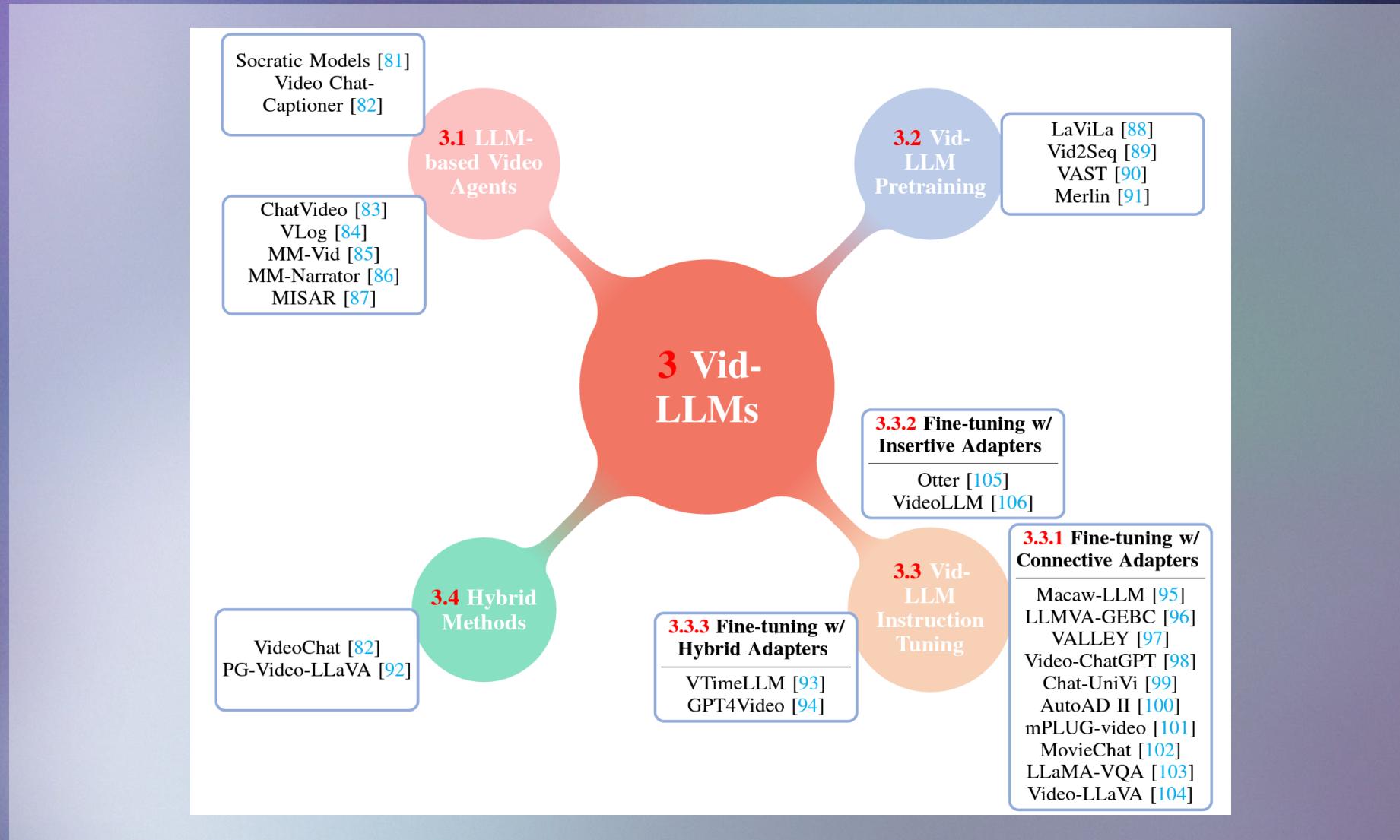
# Video Models Development



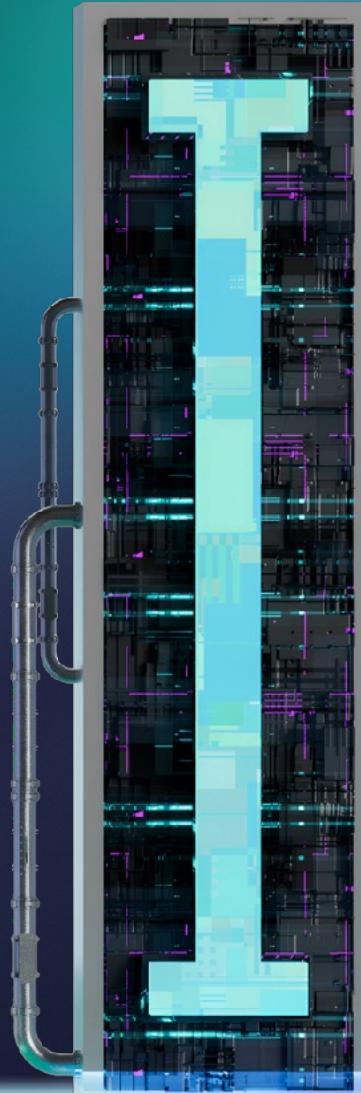
# Video-LLMs Evolution



# Video-LLMs Taxonomy



# Emotional Fusion Brain Challenge 4.0



# Fusion Brain Challenge Evolution

	Fusion Brain Challenge 1.0 AI Journey 2021	Fusion Brain Challenge 2.0 AI Journey 2022	Fusion Brain Challenge 3.0 AI Journey 2023	Emotional Fusion Brain Challenge 4.0 AI Journey 2024
<b>Modalities</b>	2	2 Images Text	3 Images Text & Audio	3+ Video (Images) Text & Audio
<b>Open Tasks</b>	4	6	Multiple	3
<b>Hidden Tasks</b>	-	6	-	-
<b>Features</b>	All tasks on natural language  Has hidden tasks	All tasks on natural language  Not limited number of tasks	All tasks on natural language  Social interactions	Has hidden metric

# Task

To develop a universal multimodal model to handle **three input modalities: video, audio and text**.  
Contest focuses on videos containing recordings of **social interactions and human emotions**.

## Task types

### Video Question Answering

This task requires answering a question with a multiple-choice answer based on the video content.



### Audio-Video Question Answering

This task requires answering a question based on the video and audio (speech/human sounds/nature) content.



### Video Captioning

This task requires a detailed description of the video, including human emotions and interactions.



# Metrics

Two metrics: **classification (Accuracy) for QA tasks and generative (METEOR) for Captioning task.**

- **Accuracy metric** (answer accuracy rate) - to evaluate the quality of answers to multiple-choice questions. It is based on an internal assessment of the model's confidence in each of the answers.
- **METEOR metric** - for evaluating the model response generation. It is based on the analysis of n-grams and focused on the use of a statistical and accurate evaluation of the source text.

Those metrics will be integrated into final one – **Integral metric**.

$$I = w_M \times \sum_j^{J_0} METEOR_j + w_A \times \sum_i^{J_m} Accuracy_i$$

where  $J_0$  and  $J_m$  - the number of questions of each type,  
 $w_M$  and  $w_A$  are the weights of the METEOR and Accuracy metrics, respectively.

# Input Data Format

Each data object is represented as a dictionary with the following fields defined:

- `task_id` - (type: int) represents a unique identifier for the question;
- `task_type` - (type: str) represents the type of task, which can be one of the following options: `captioning` / `qa`;
- `question` - (type: str) contains the question related to the video;
- `video` - (type: str) describes the path to the video file;
- `audio` - (type: str) describes the path to the audio file, not filled in all tasks;
- `choices` - (type: list) this field contains a list of answer options for the given question.

Each answer option:

```
{'choice_id': 0, 'choice': 'value'},
```

where `'choice_id'` is the sequential index of the answer option from range [0, 4],  
and `'choice'` is the actual value of the answer option.

Videos for which predictions need to be made are placed in the ``video`` folder in `.mp4` format.  
Audio recordings are placed in the ``audio`` folder in `.mp3` format.

# Input Data Format

Audio-Video QA

```
{ 'task_id': 1,  
  'task_type': 'qa',  
  'question': 'Who is dancing in the end of the video?',  
  'video': 'path_to_video.mp4',  
  'audio': 'path_to_audio.mp3',  
  'choices': [  
      {'choice_id': 0, 'choice': 'Woman in red'},  
      {'choice_id': 1, 'choice': 'Woman in blue'},  
      {'choice_id': 2, 'choice': 'Man in green'},  
      {'choice_id': 3, 'choice': 'Man in black'},  
      {'choice_id': 4, 'choice': 'Nobody'}  
  ] }
```

Video Captioning

```
{ 'task_id': 2,  
  'task_type': 'captioning',  
  'question': 'Describe this video in detail.',  
  'video': 'path_to_video.mp4',  
  'audio': '',  
  'choices': [] }
```

# Submission file structure

Videos for which predictions need to be made are placed in the `video` folder in **.mp4** format.  
Audio recordings are placed in the `audio` folder in **.mp3** format.

```
└── run.py
    ├── video/
    ├── audio/
    ├── dataset_file.json
    └── README.md
    └── your-code-folder
        ├── evaluate.py
        ├── mm_utils.py
        ├── constants.py
        └── ...
```

# Output Data Format

The output files should be saved to `output\_path\_from\_job` directory (automatically created in the root of the solution).

- ./output\_path\_from\_job/for\_generative\_metric.json file for **Video Captioning** task ;
- ./output\_path\_from\_job/for\_classification\_metric.json for **Question Answering**.

The answer values are the index of the correct answer from the provided answer choices, starting with 0 and ending with index 4;

## Video Captioning

```
{  
  0: "Video shows ...",  
  3: "Here we see ..."}
```

## Video(-Audio) QA

```
{  
  1: 1,  
  2: 4,  
  4: 0}
```

# Submission Technical Details

- Default image cannot be changed.
- Dockerfile and requirements for it are on Data tab of [Official Contest page](#).
- Additional requirements can be installed through addition of libraries source code to submission.
- Additional model weights (if needed) should be included in user submission.
- The overall **user submission size should be not more than 5 Gb** (due to technical limits).
- The **disk limit is 10 Gb** (including 2.5 Gb for test dataset).
- In Docker image we stored a list of popular VidLLMs parts. The complete list of them and paths can be also found on Data tab of [Official Contest page](#) (PATHS.md).
- Resources: 243 Gb RAM, 16 CPU-cores, 1 GPU Tesla A100 (80 Гб). Max. evaluation time: 3,5 hours

## Data

 Dockerfile  
Docker-образ

 PATHS.md  
Файл с путями к весами моделей в...

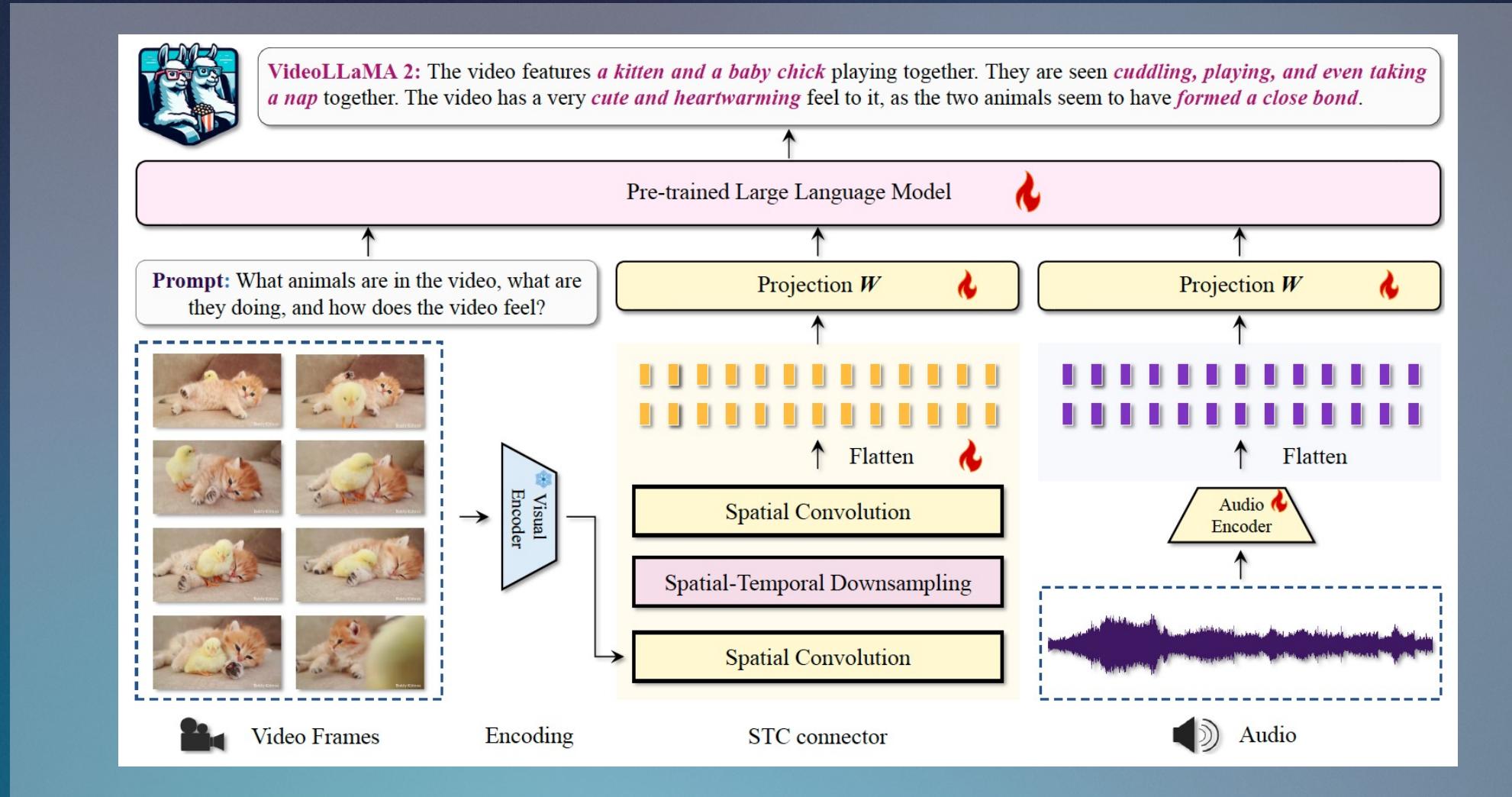
 requirements.txt  
Список установленных библиотек в...

 baseline\_submission.zip  
Baseline submission

# Model Options

Model name	Path in image	Link to HF
DAMO-NLP-SG/VideoLLaMA2-7B	/app/DAMO-NLP-SG/VideoLLaMA2-7B	<a href="#">link</a>
imagebind_huge.pth	/app/imagebind_huge.pth	<a href="#">link</a>
InternViT-300M-448px	/app/InternViT-300M-448px	<a href="#">link</a>
internlm2-chat-1_8b	/app/internlm2-chat-1_8b	<a href="#">link</a>
LanguageBind_Audio_FT	/app/LanguageBind_Audio_FT	<a href="#">link</a>
llava-onevision-qwen2-0_5b-sis	/app/llava-onevision-qwen2-0_5b-sis	<a href="#">link</a>
Qwen2-0.5B-Instruct	/app/Qwen2-0.5B-Instruct	<a href="#">link</a>
llava-onevision-qwen2-7b-sis	/app/llava-onevision-qwen2-7b-sis	<a href="#">link</a>
clip-vit-large-patch14-336	/app/clip-vit-large-patch14-336	<a href="#">link</a>

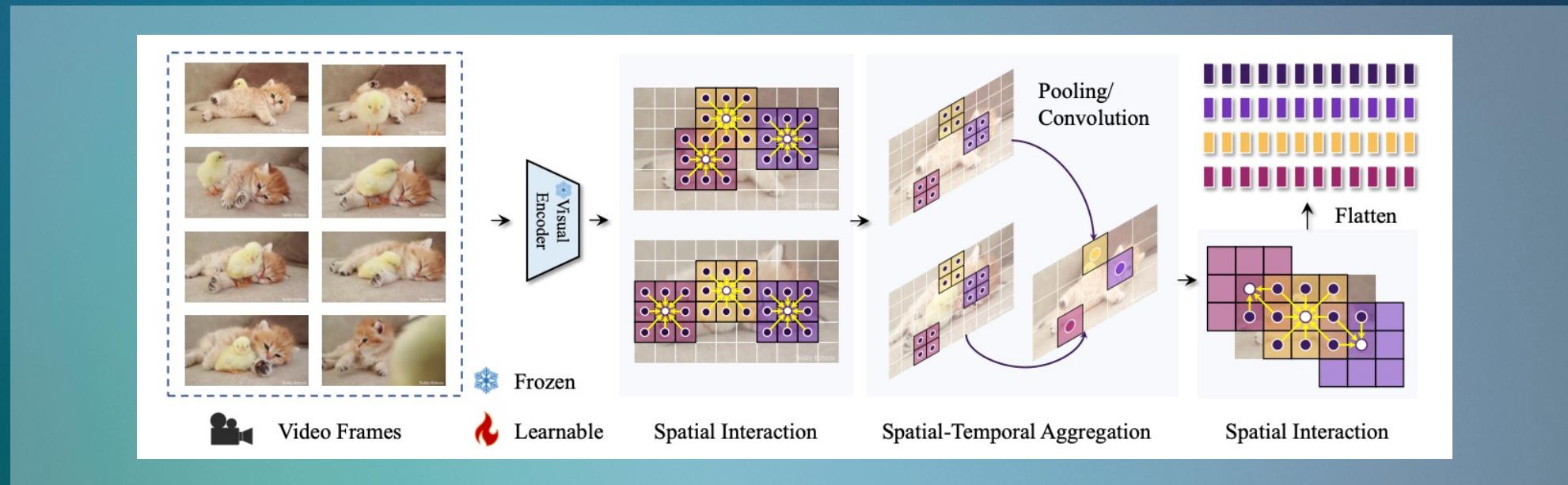
# Baseline Solution



# Baseline Solution

## Spatio-temporal video encoding. STC Connector

- STC consists of **two** spatial interaction modules (RegStage adaptation) and **one** spatial-temporal aggregation module (3D convolution adaptation)
- Maintain the spatial-temporal order in the output visual tokens
- Reduce the number of spatial-temporal tokens
- Alleviate information loss during spatial-temporal downsampling



# Additional nominations

## The Role Game\*

- to apply the developed multimodal model to **the video recordings of the role game**
- the model will have to answer questions about:
  - the course of the game,
  - determine the participants' roles
  - the likelihood of plausibility of the points voiced and demonstrated by the players.
- the Accuracy classification metric

## The Fastest Solution\*

- the models will be evaluated based on the lowest-time inference metric  $T$

$$T = \frac{1}{N} \sum_{i=1}^N \tau_i$$

where  $\tau_i$  is the inference time for task  $i$ , sec.  
 $N$  is the number of tasks.