

Research Project Proposal:

3D TinyML for Video Streaming Analysis

Hazem Shalby
Hazemhesham.shalby@mail.polimi.it
CCS Track



POLITECNICO
MILANO 1863



HP-SR
in Information Technology

3D TinyML for Video Streaming Analysis

How is it possible to design 3D Tiny Machine Learning solutions able to support a video streaming analysis on tiny devices technologically constrained on memory, computation, and energy?

3D TinyML for Video Streaming Analysis

Tiny machine learning is the field of **machine learning** technologies including hardware, algorithms and software capable of performing on-device analytics at extremely low power (\approx mW), enabling a variety of always-on use-cases and targeting battery operated devices.

- ✓ **Increase autonomy**

- ✓ **Reduce decision-making latency**

- ✓ **Reduce transmission bandwidth**

- ✓ **Increase energy-efficiency**

- ✓ **Security and Privacy**

- ✓ **Incremental/Adaptive Learning**

- ✓ **Ecosystem of units**

- **Low computing ability**

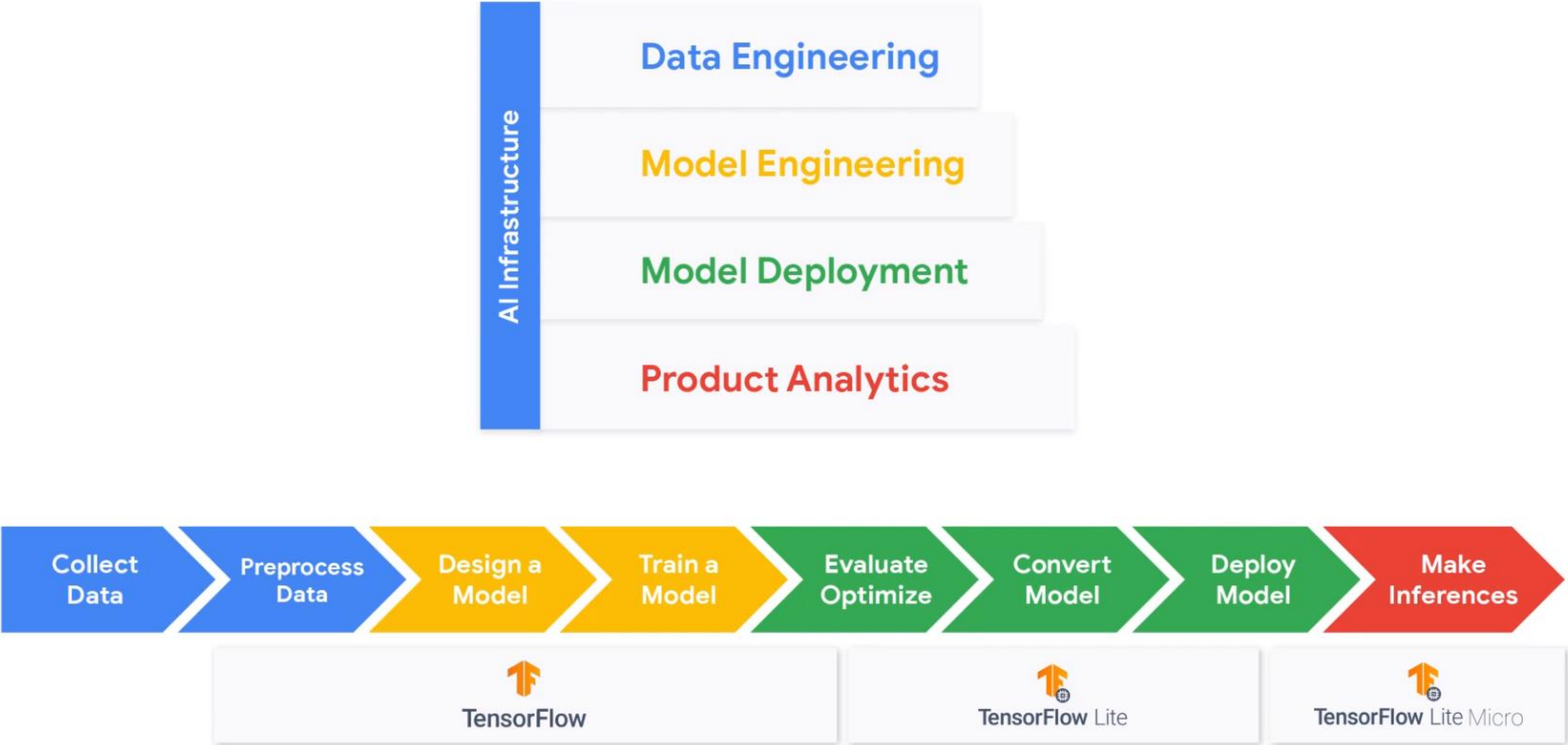
- **Constraints on energy**

- **Constraints on memory (RAM/FLASH)**

- **Complexity in design and development**

- **Strong connection between HW, SW and ML**

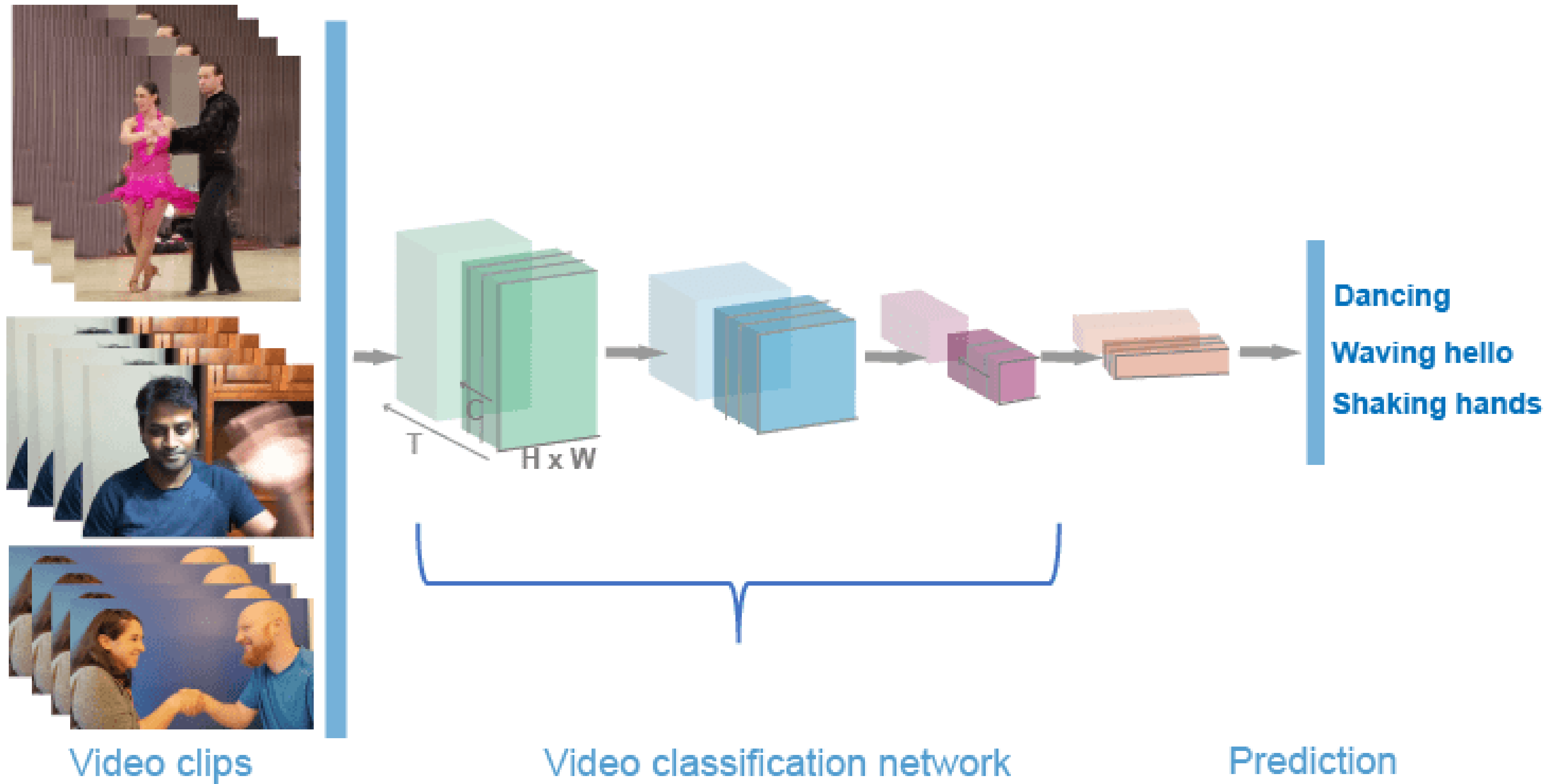
3D TinyML for Video Streaming Analysis



3D TinyML for Video Streaming Analysis

Being $x_t \in U \subset R^{M \times N \times C}$ (M, N, and C are the sizes of the input) a frame of the video streaming and $y_t \in Y = \{\Omega_1, \dots, \Omega_c\}$ a label associated to x_t , the goal is to construct a classifier $f_\theta(x_t, x_{t-1}, x_{t-2}, \dots)$ able to map an unseen data \bar{x}_i to its label \bar{y}_i .

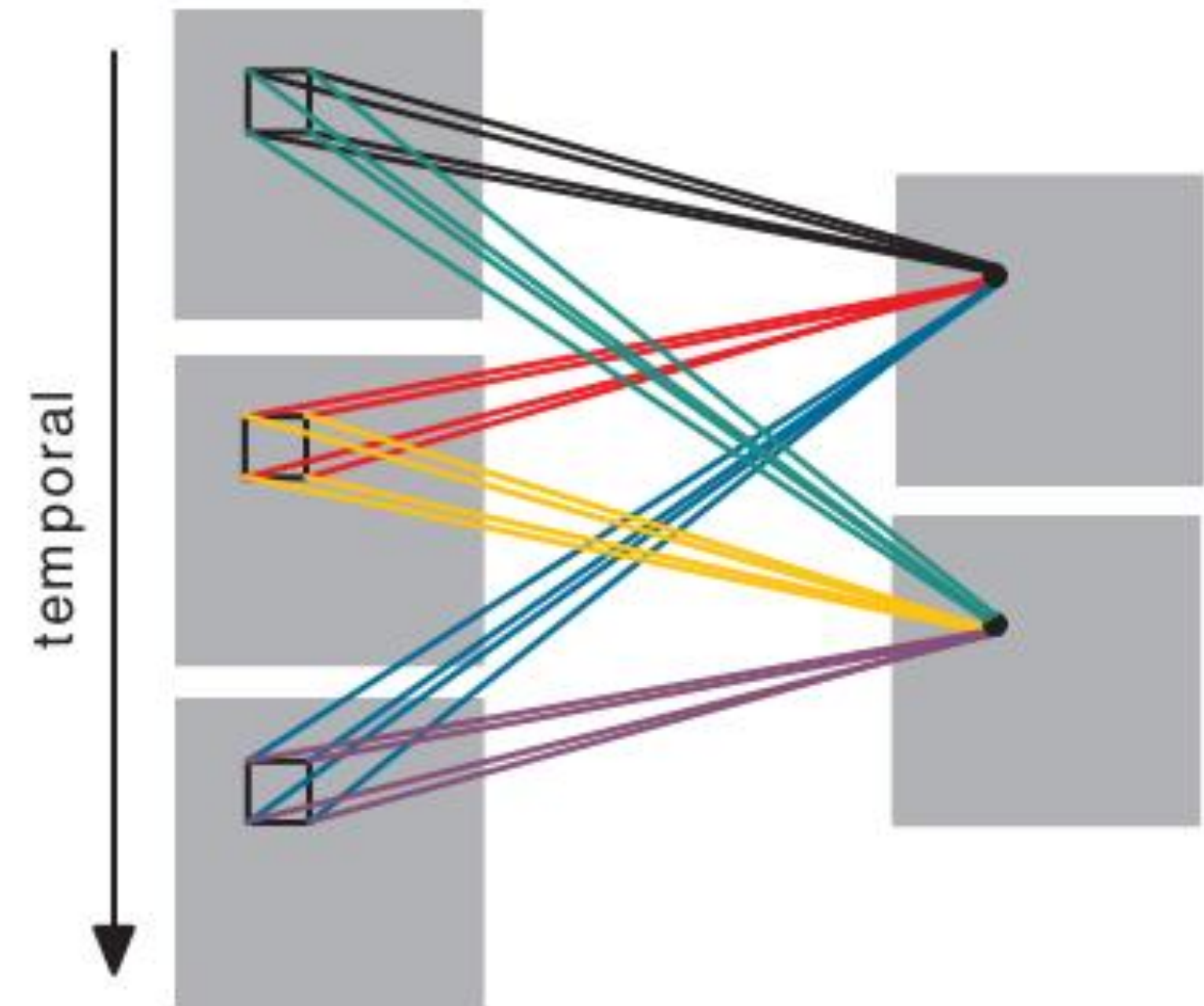
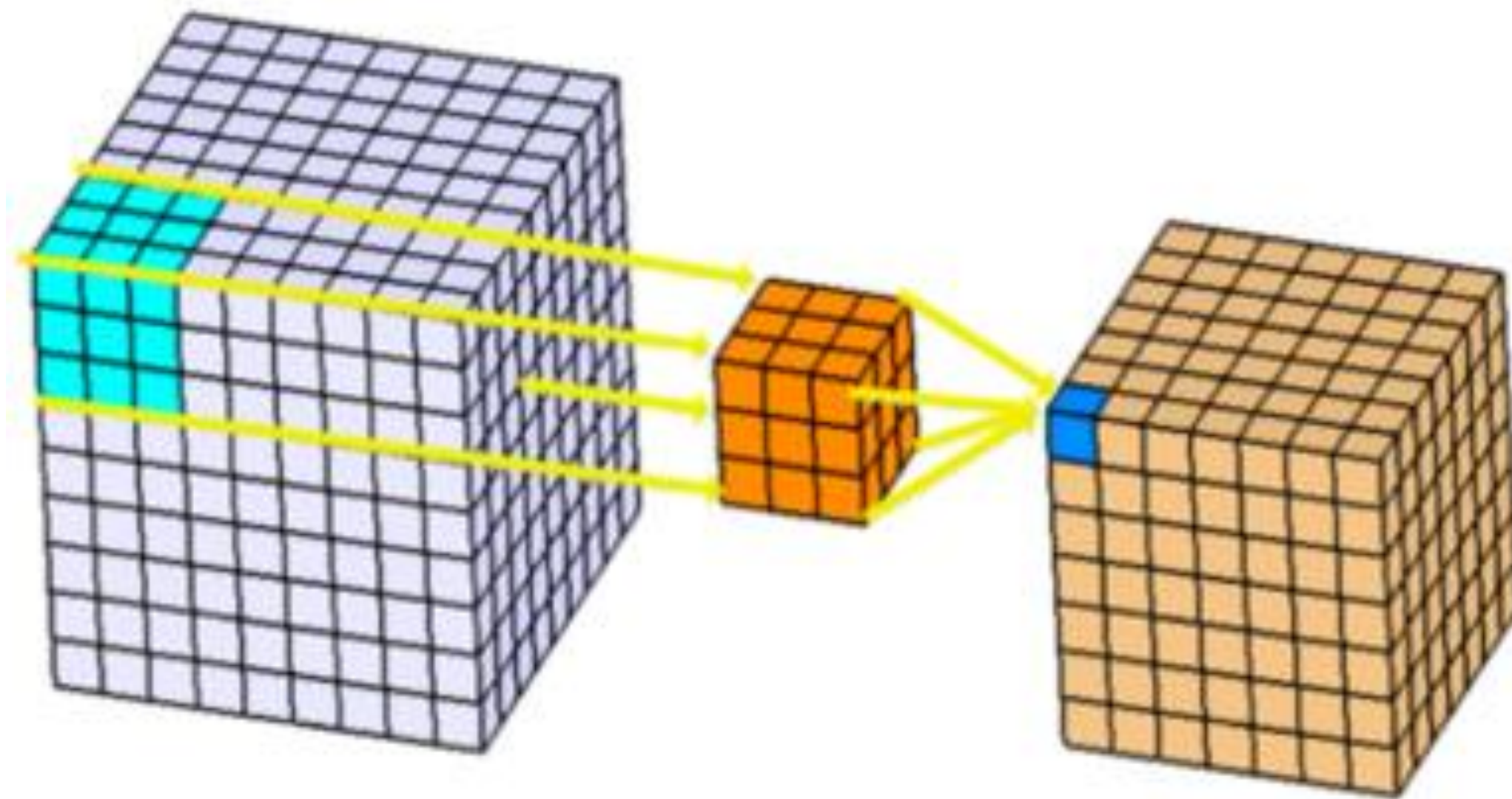
3D TinyML for Video Streaming Analysis



3D TinyML for Video Streaming Analysis

3D convolutions applies a 3-dimensional filter to the dataset and the filter moves 3-direction (x, y, z) to calculate the low-level feature representations.

They are helpful in **event detection in videos**.

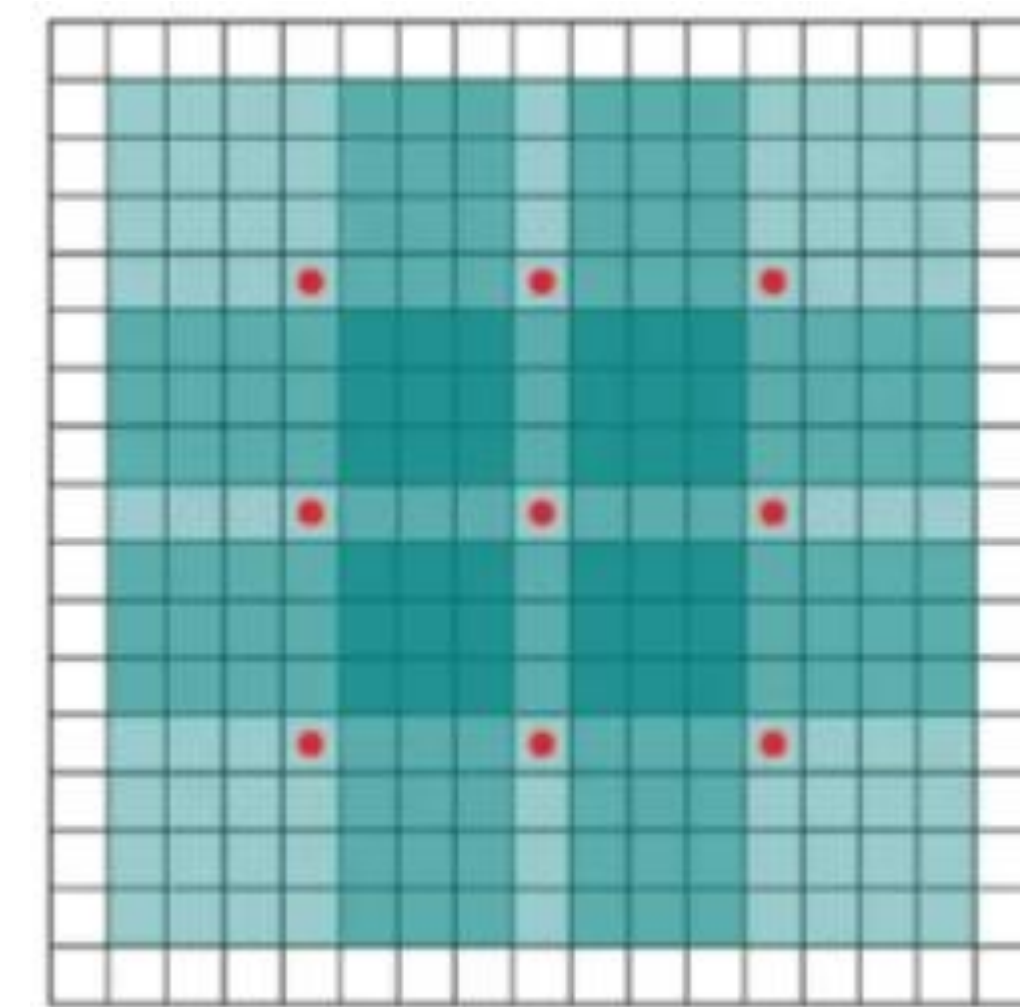
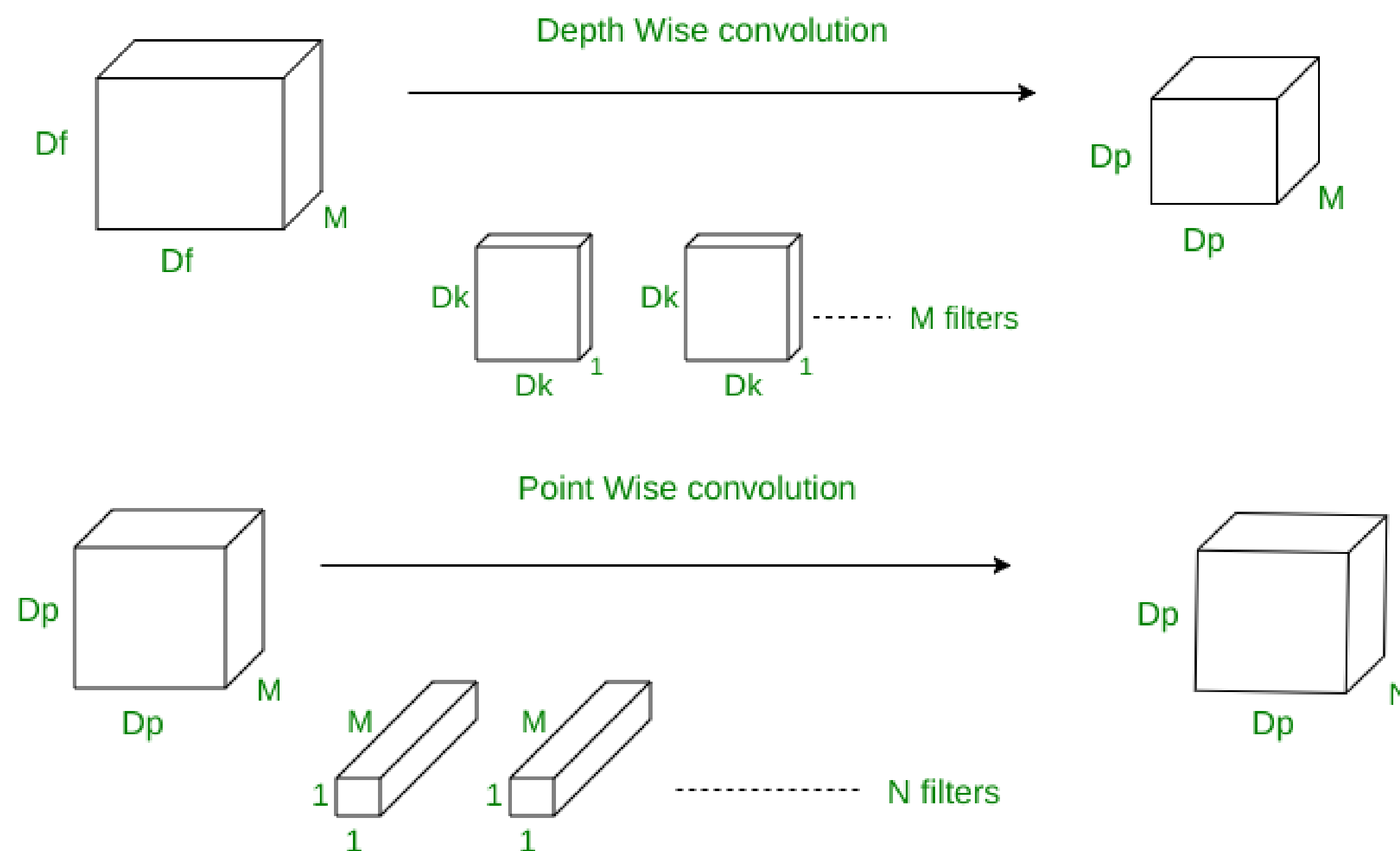


Research Steps

- Redesigning the CNN architecture
- Introducing approximate computing mechanisms
- Exploiting embedded-system code optimization

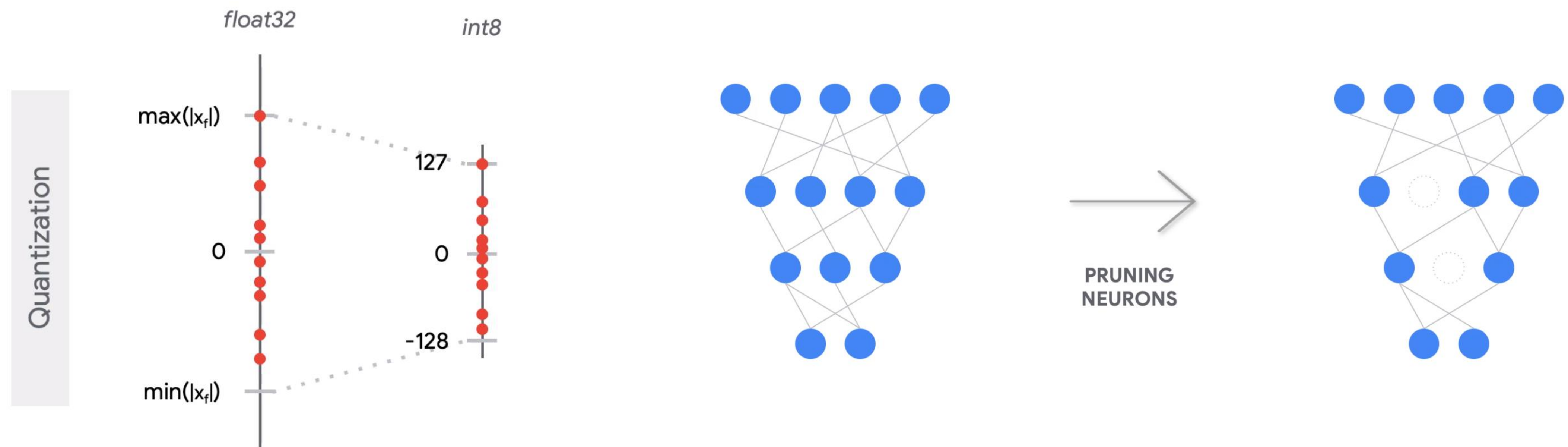
Research Steps: CNN architecture

Explore existing and efficient ways of doing the standard 2D convolutions (e.g. **separable convolutions** and **dilated convolutions**) and combine them with **3D convolutions**, which are suitable for the multiple frame analysis task but implementing them in the TinyML field is challenging.



Research Steps: approximation

- **Precision scaling:** reduce the memory occupation of a CNN by changing the precision
- **Task dropping:** reduce the computational load and memory occupation by skipping the execution of certain tasks associated with the processing pipeline.

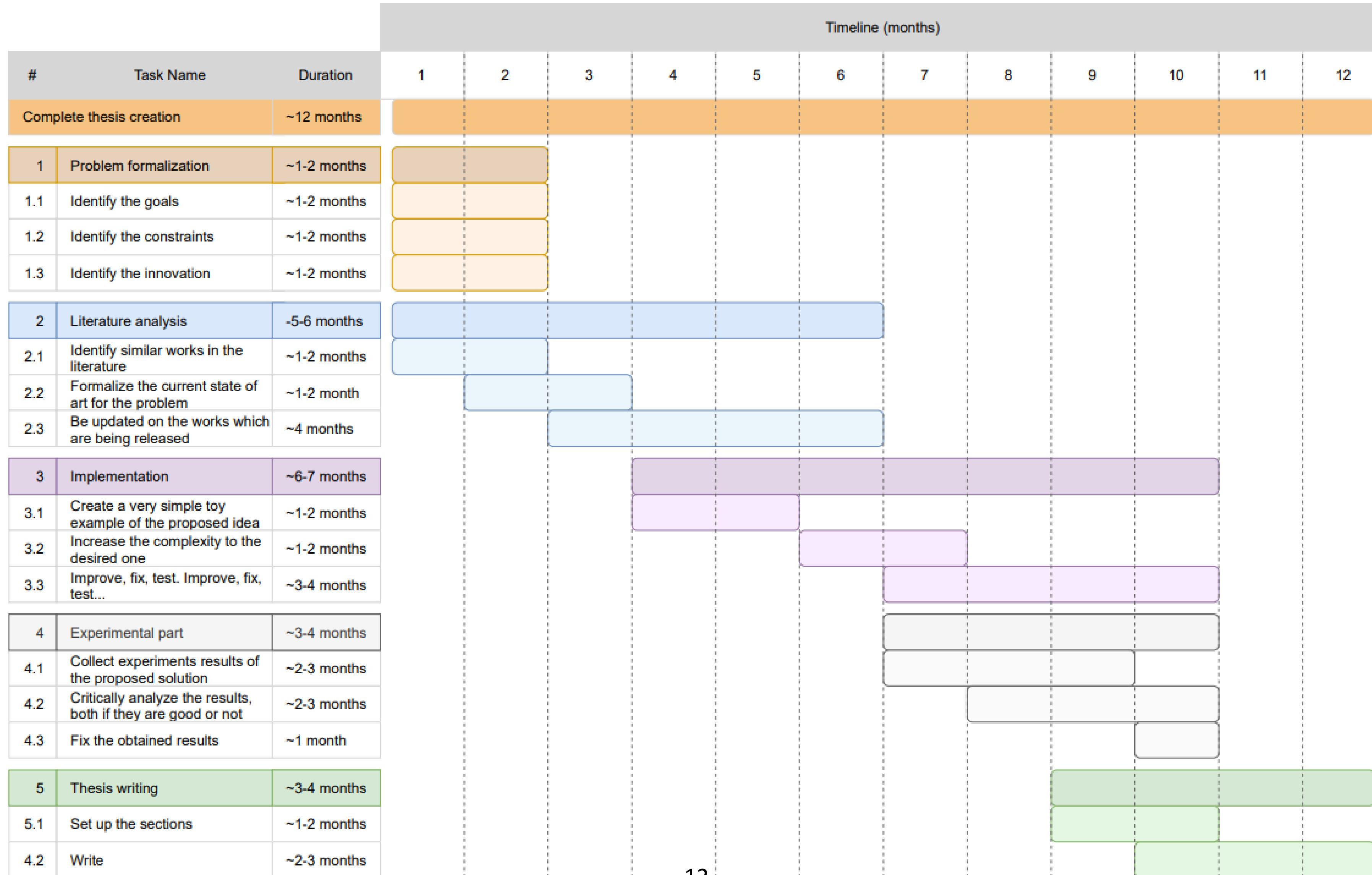


Research Steps: Code Optimization

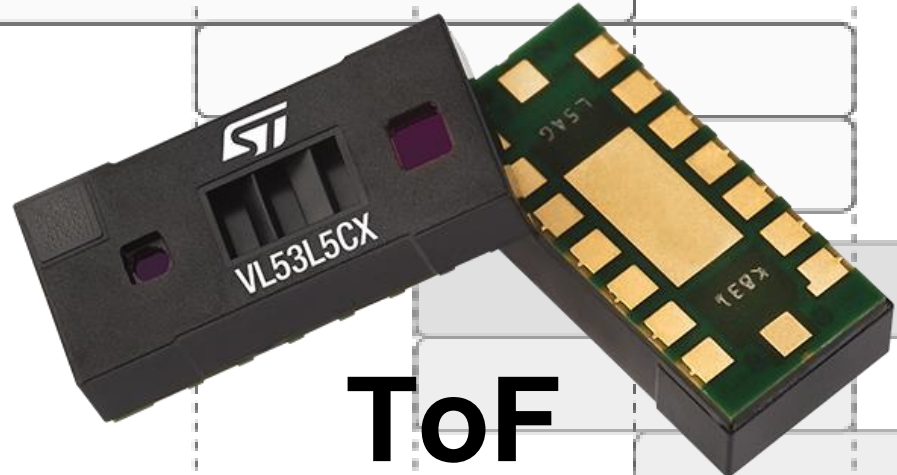
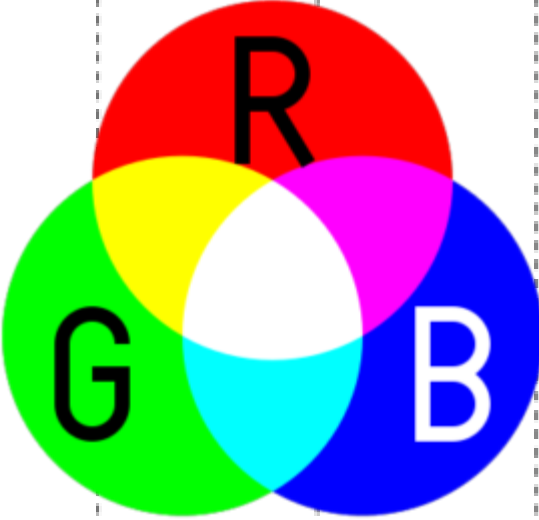
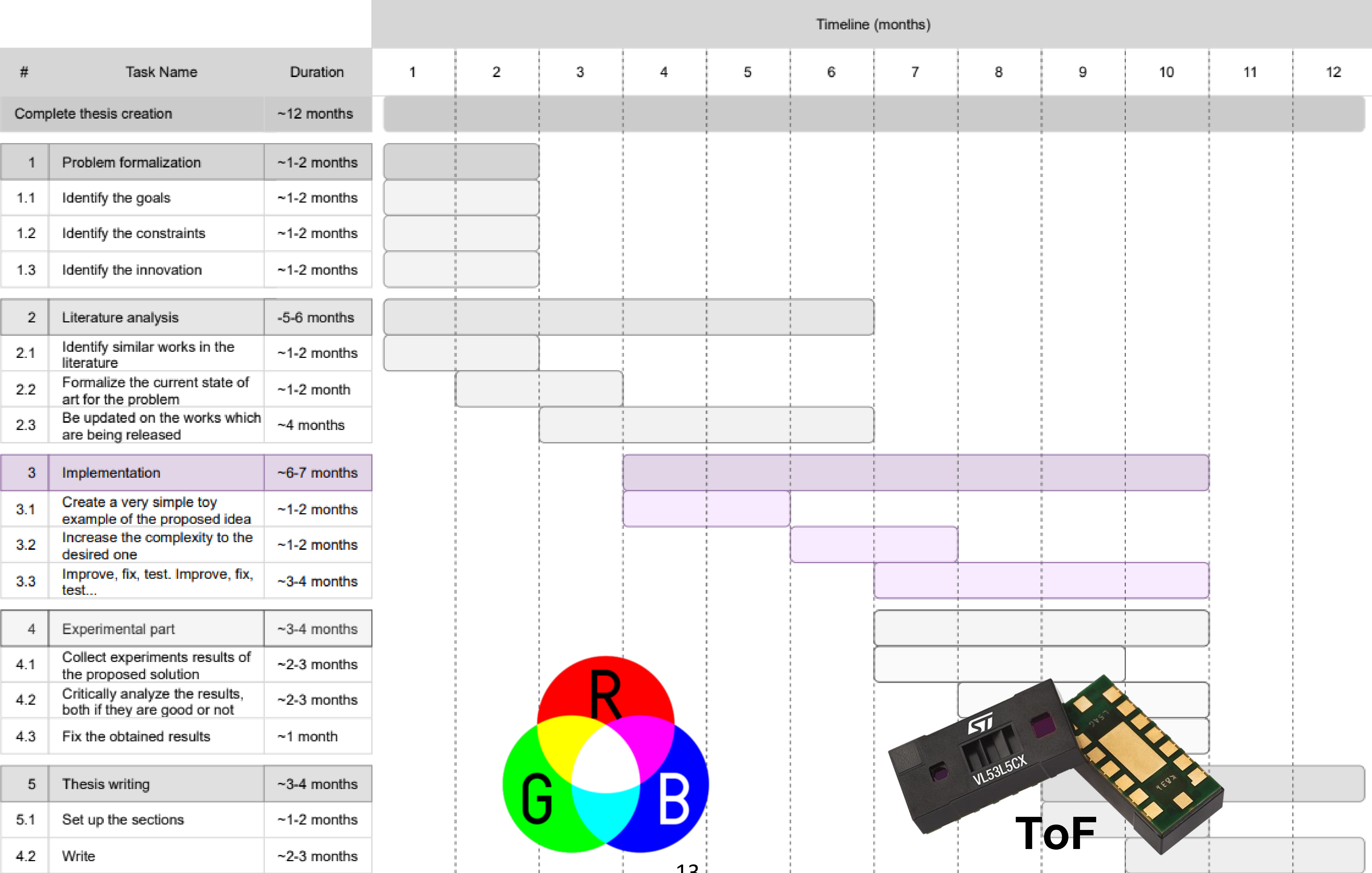
involves the use of **toolchains** or **code optimization** mechanisms to further reduce the computational and memory demands for a target hardware platform



Research Plan



Research Plan



ToF

Research Assessment

The accuracy of the classifier

The number of multiple contiguous frames used by the classifier

The time needed for the inference

Memory occupation

The number FLOPs and MACs operations needed

Thanks for your attention!