

Video Highlights Detection

Александр Гордеев
Computer Vision Engineer
SberDevices

Описание задачи

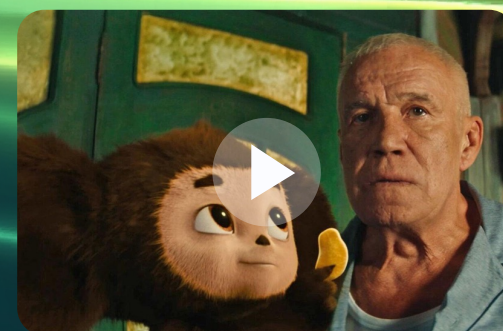
1 Полнометражный фильм



2 Black box



3 Интересные моменты



Video shorts
здесь и сейчас



Shorts

Модальности фильма



Аудио модальность



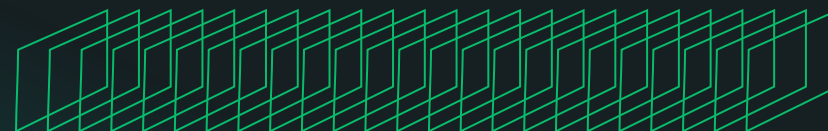
Видео модальность



Структура видео-потока

Frame

Одиночное неподвижное изображение, которое при последовательном воспроизведении с другими кадрами видеоролика создает движение на поверхности воспроизведения



Shot

Серия кадров, снятых одной и той же камерой в течение непрерывного периода времени

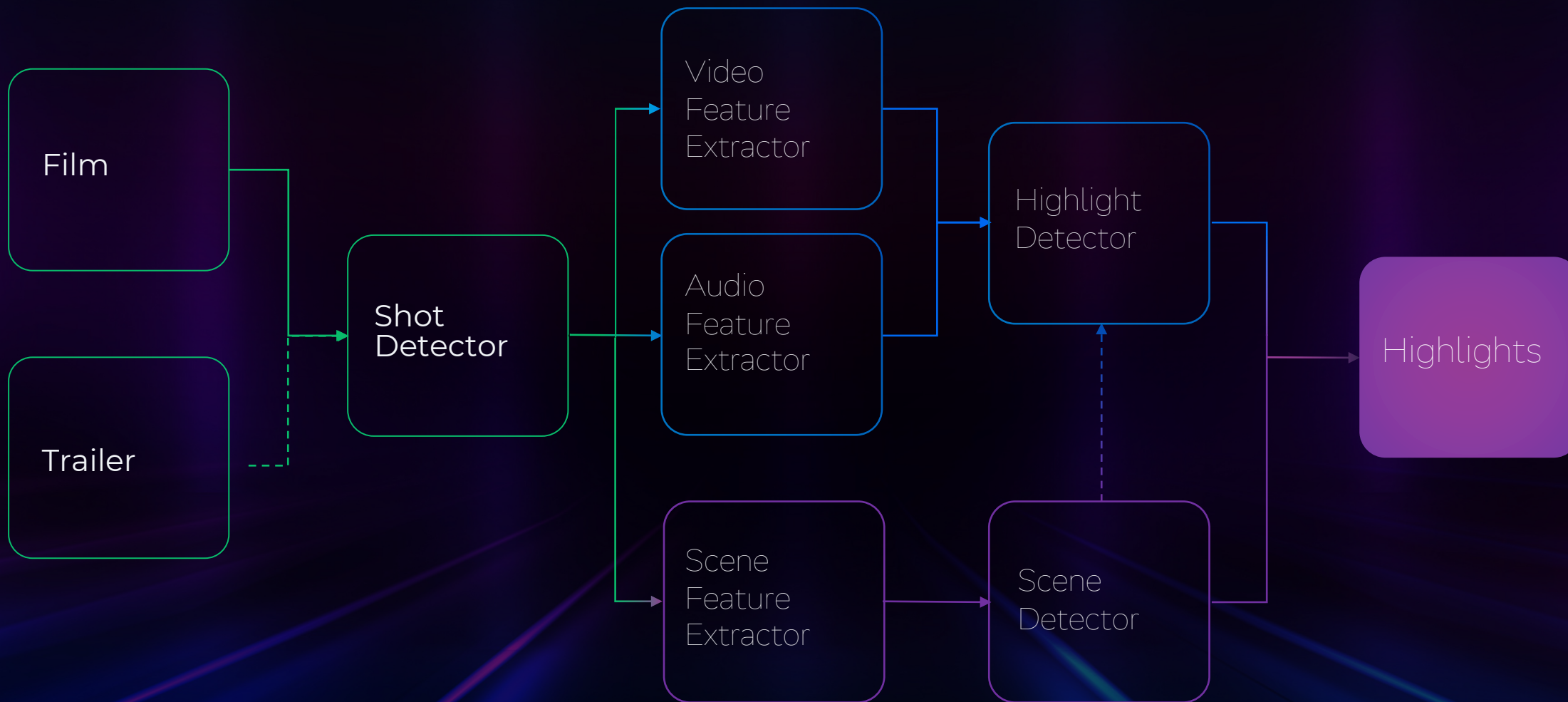


Scene

Серия shot-ов, представляющая собой семантически связную часть сюжета

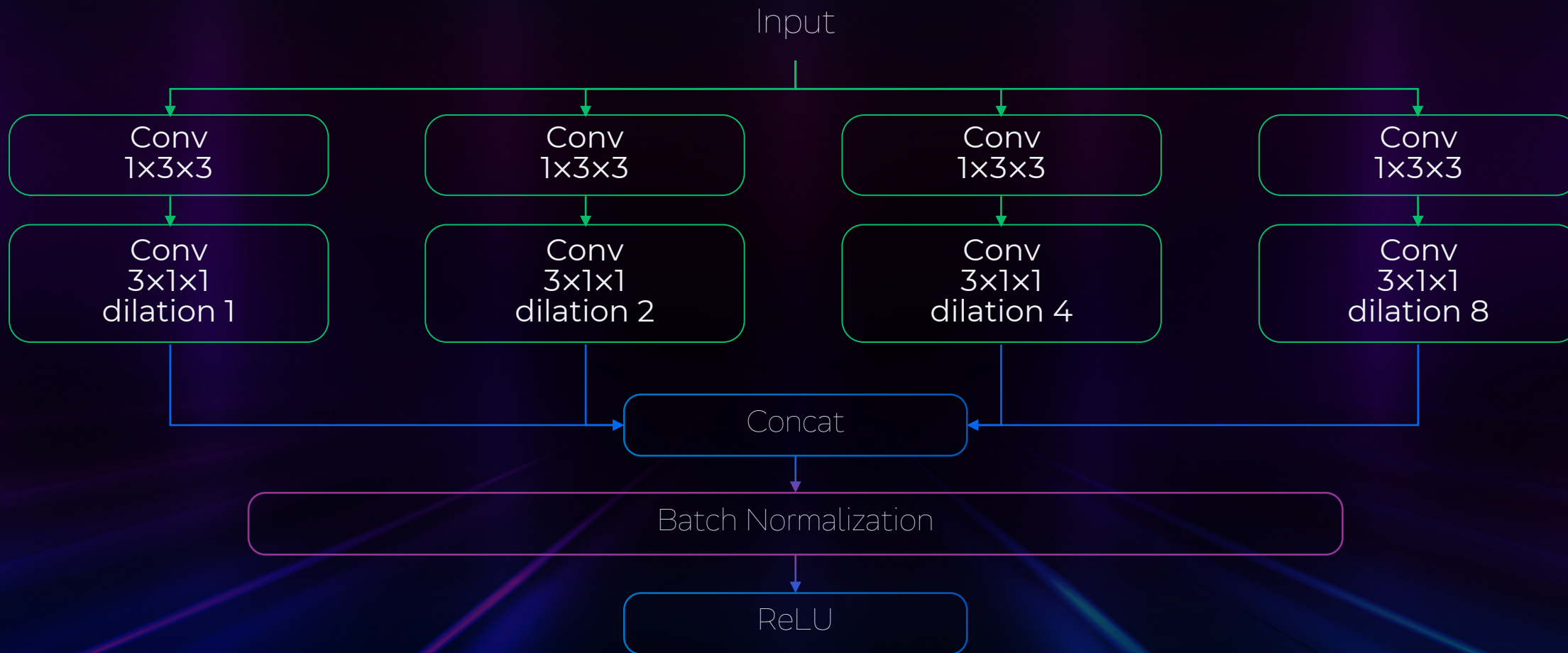


Архитектура решения



Shot Detector

Dilated 3D Conv Block

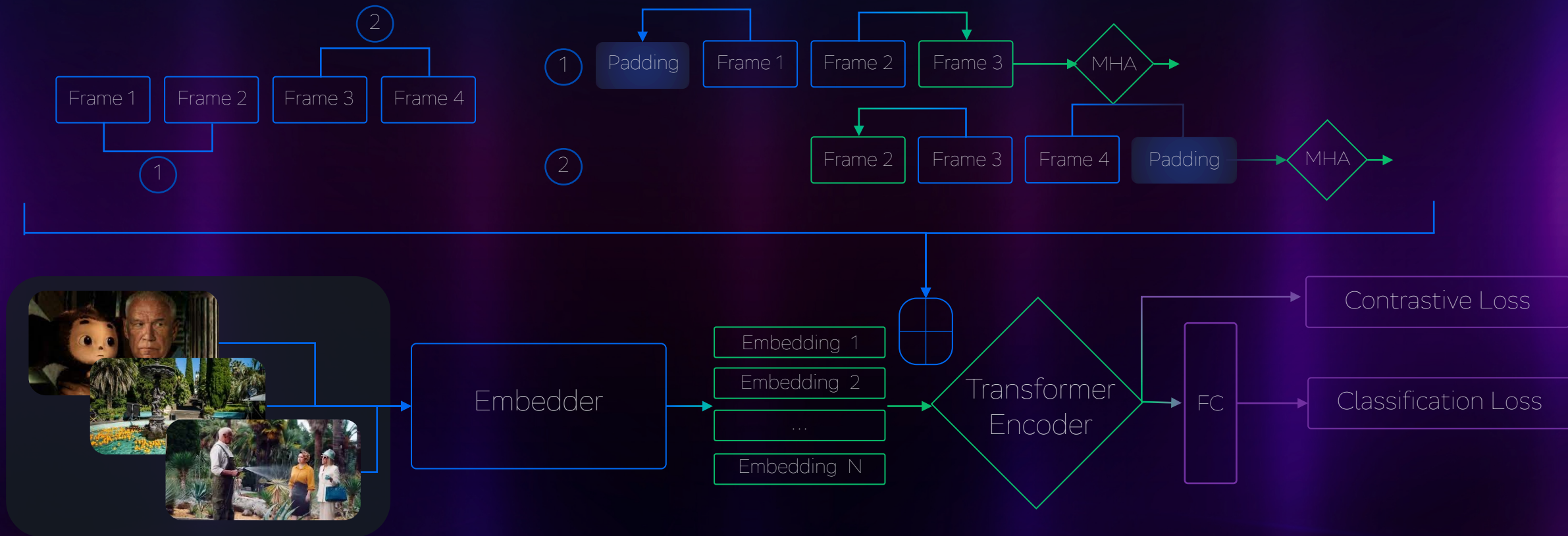


Shot Detector

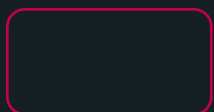
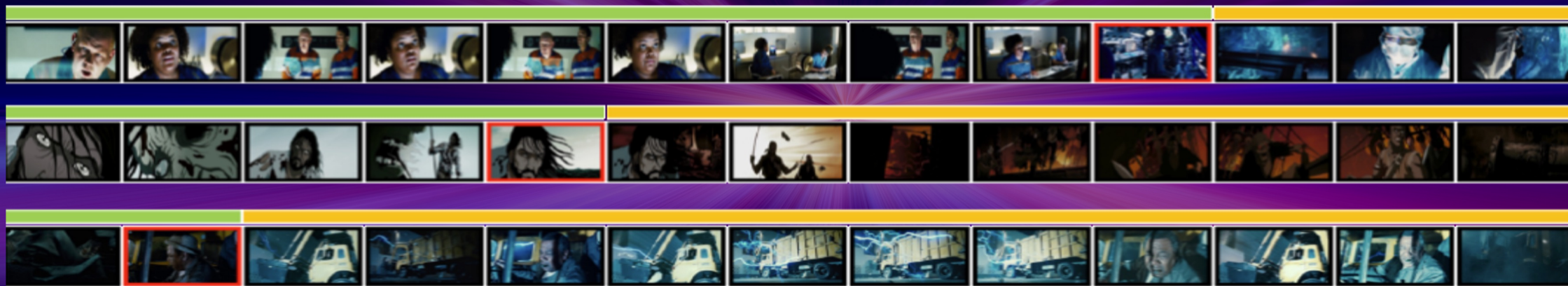
Бенчмарк:
целевая метрика —
F1-Score

Dataset	ClipShots	BBC	RAI
DSMs	0.761	0.893	0.928
ST ConvNets	0.759	0.926	0.939
TransNet	0.735	0.929	0.943
TransNetV2	0.776	0.962	0.939
AutoShot	0.787	0.971	0.955
Ours	0.808	0.980	-

Shot Detector



Scene Detector



Pseudo-boundary



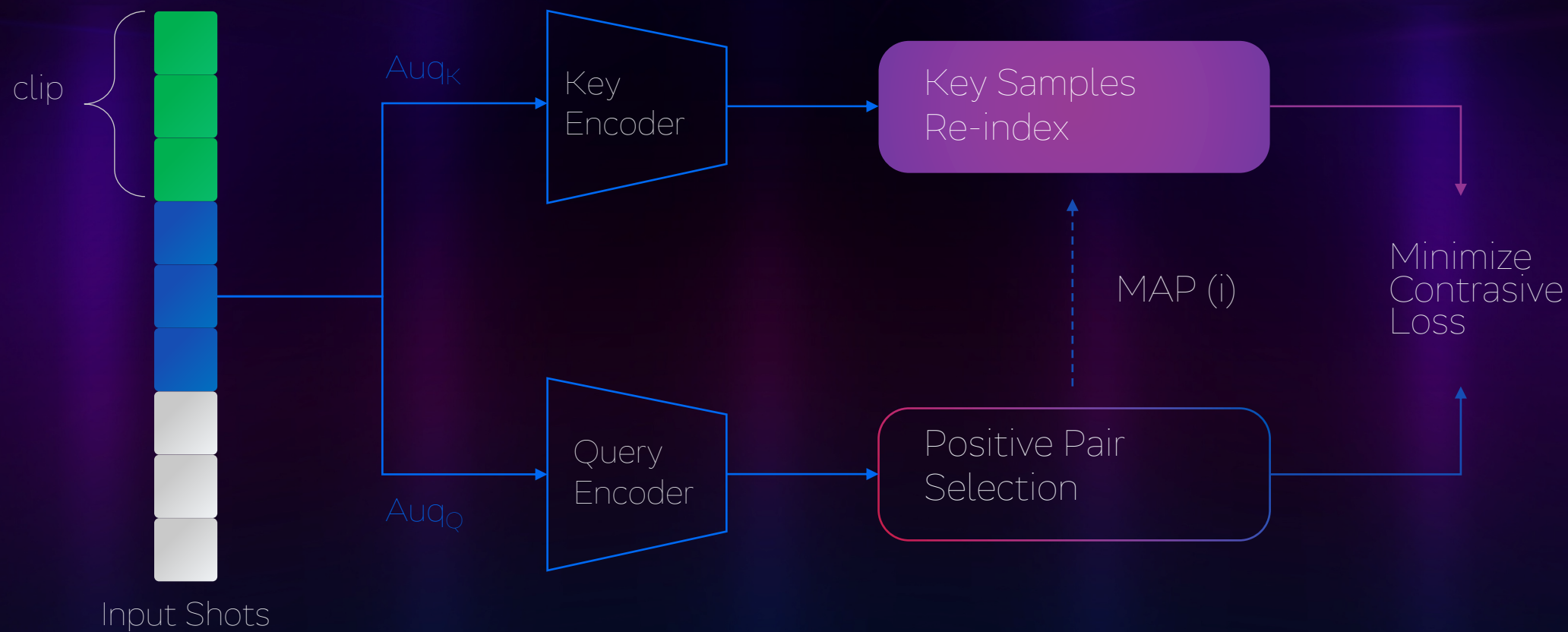
Pseudo Scene A



Pseudo Scene B

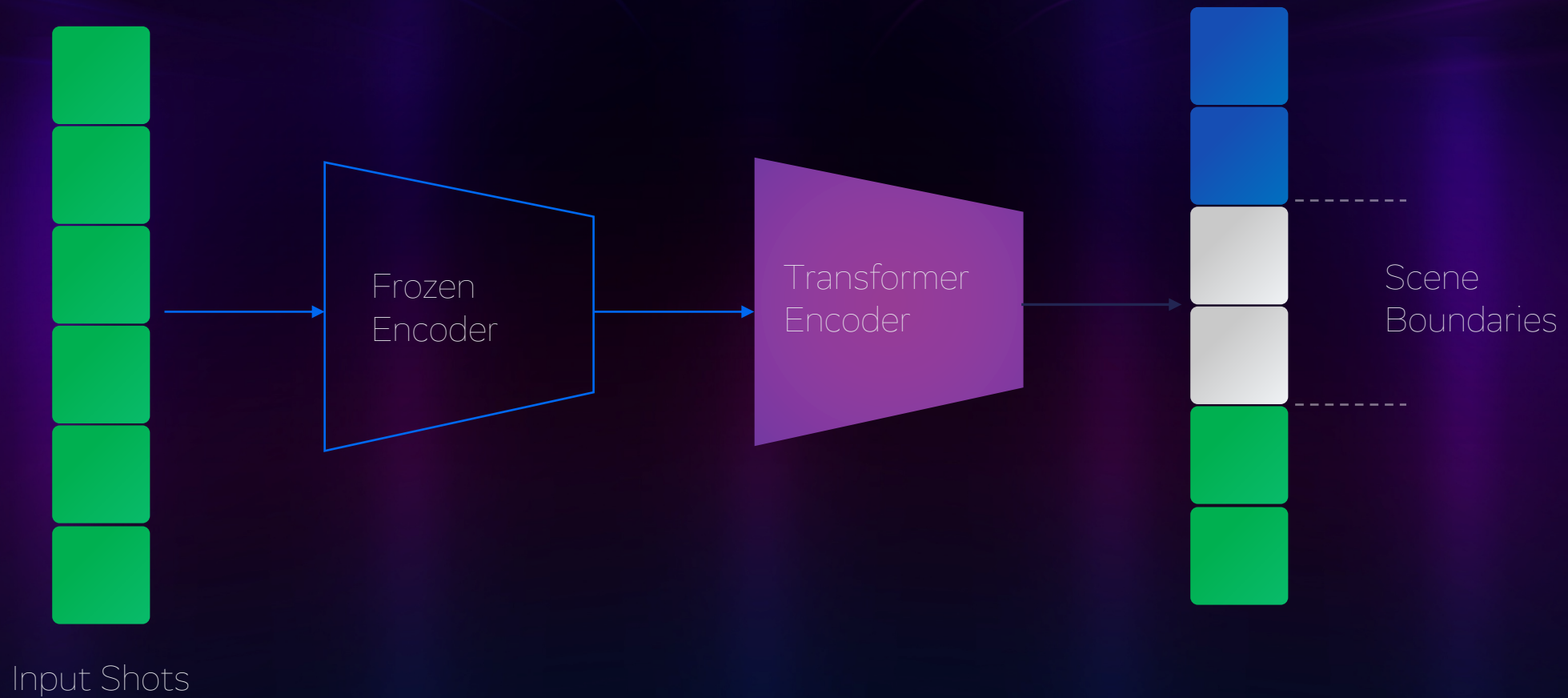
Scene Detector

Representation Learning Stage

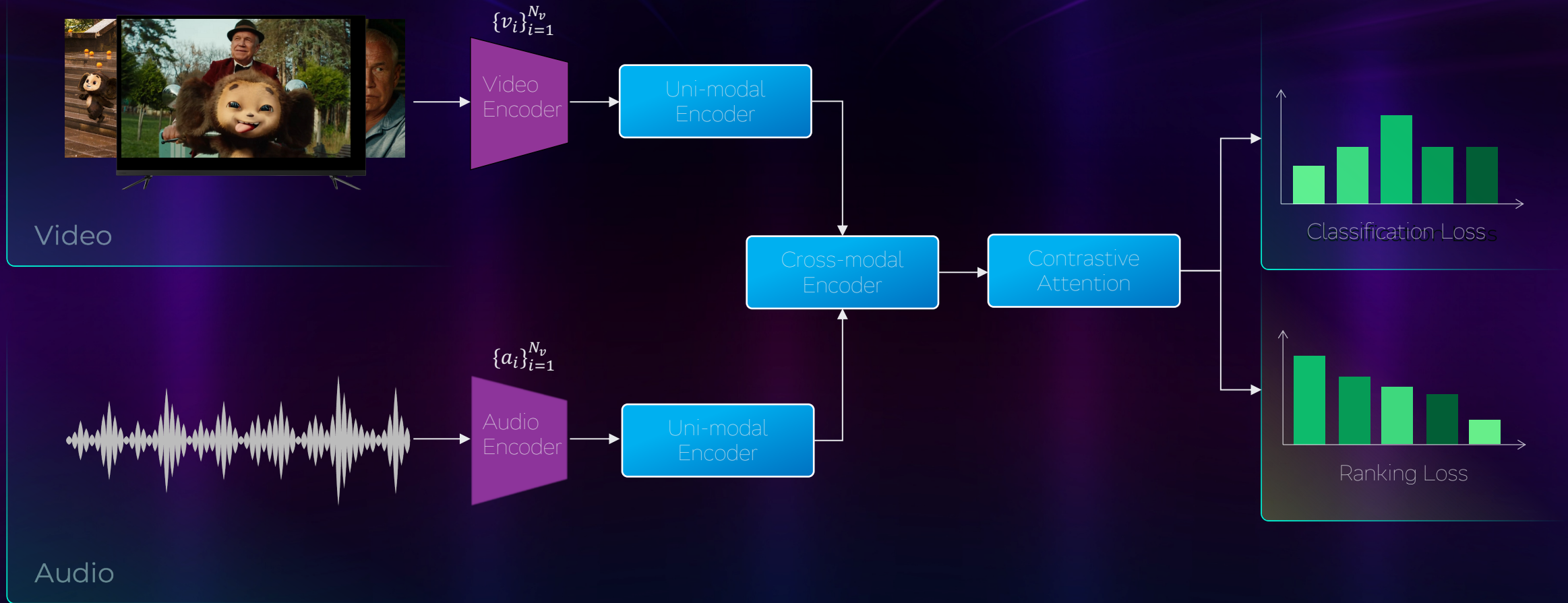


Scene Detector

Video Scene Segmentation Stage



Highlight Detector



Глобальный контекст

Local Video Embeddings



Unimodal Encoder

Global Video Embeddings



Local Audio Embeddings

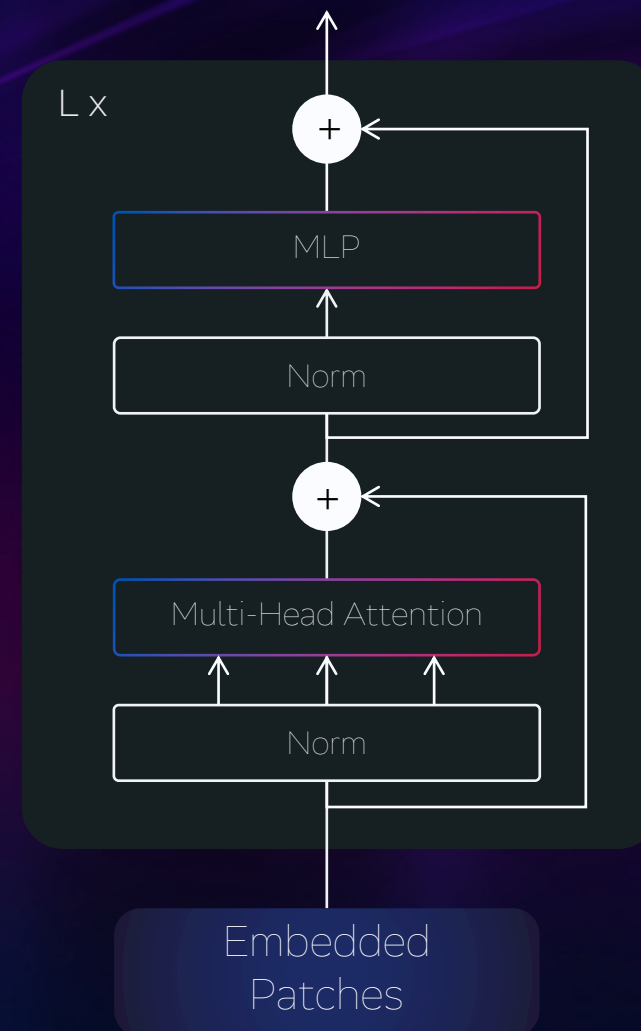


Unimodal Encoder

Global Audio Embeddings

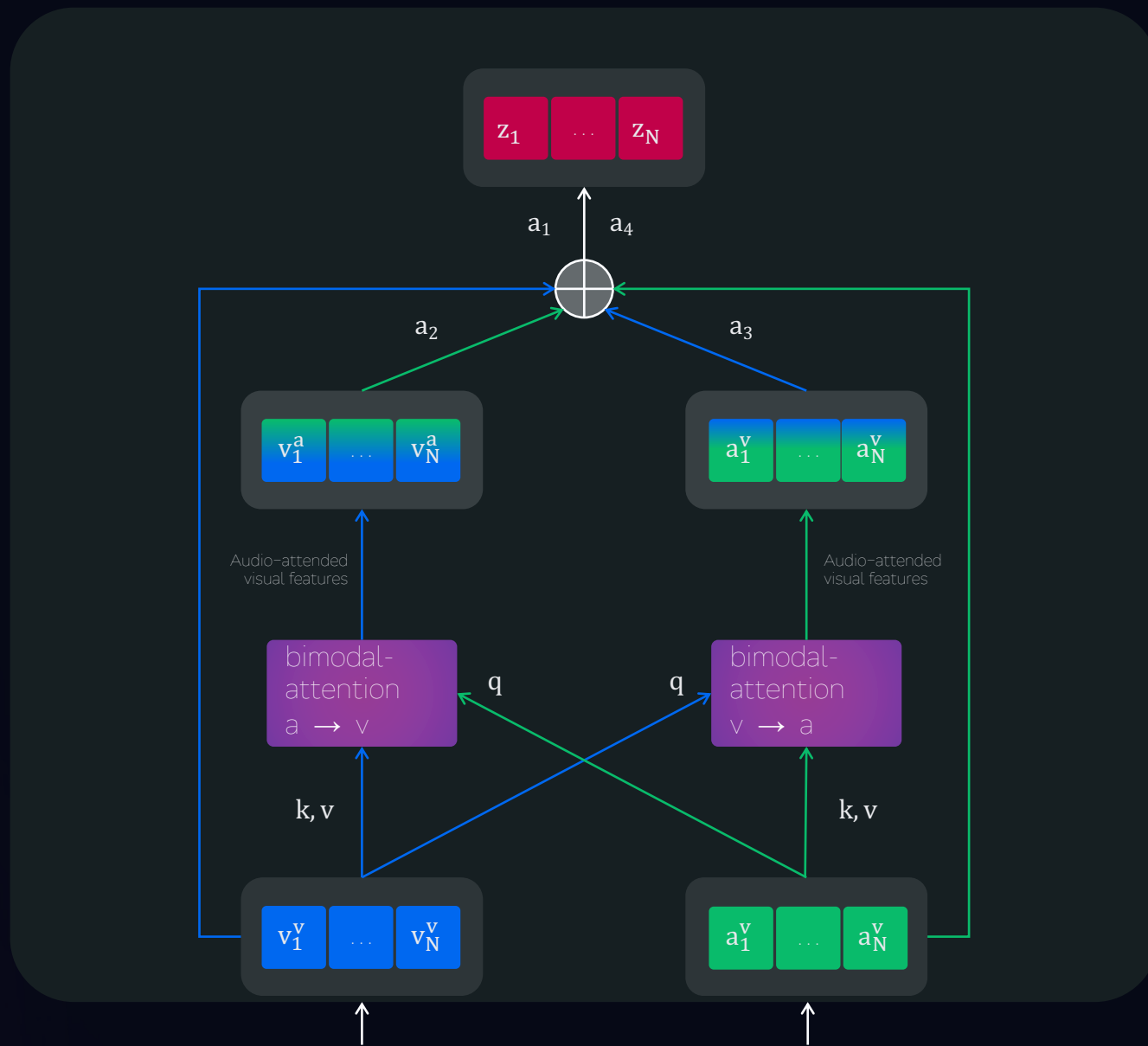


Uni-modal Encoder



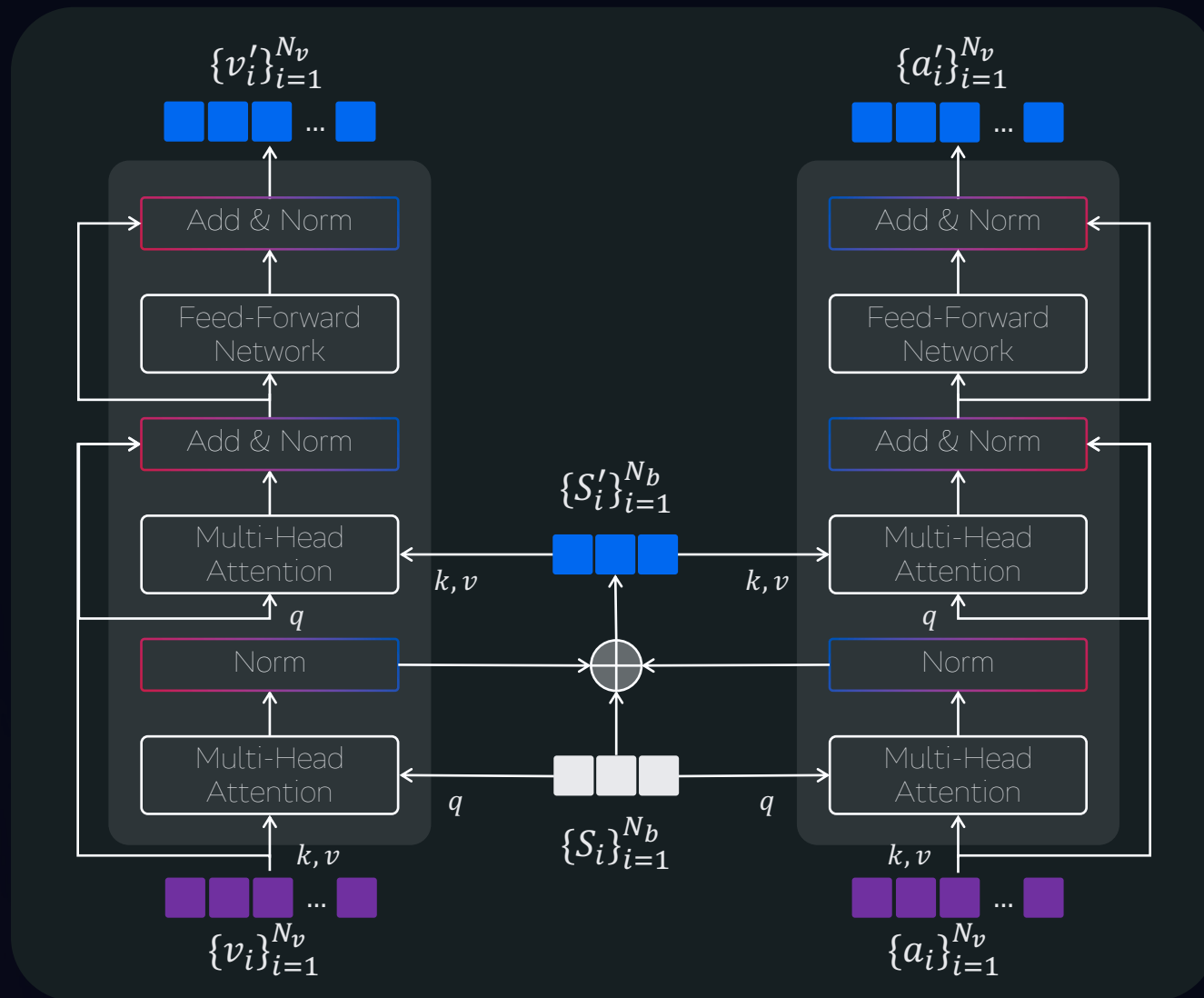
Фьюзинг модальностей

Cross-modal attention

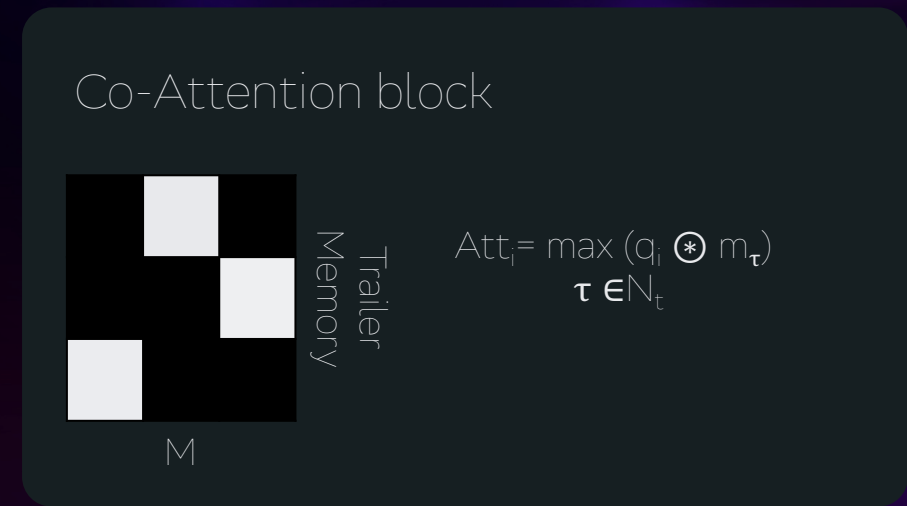
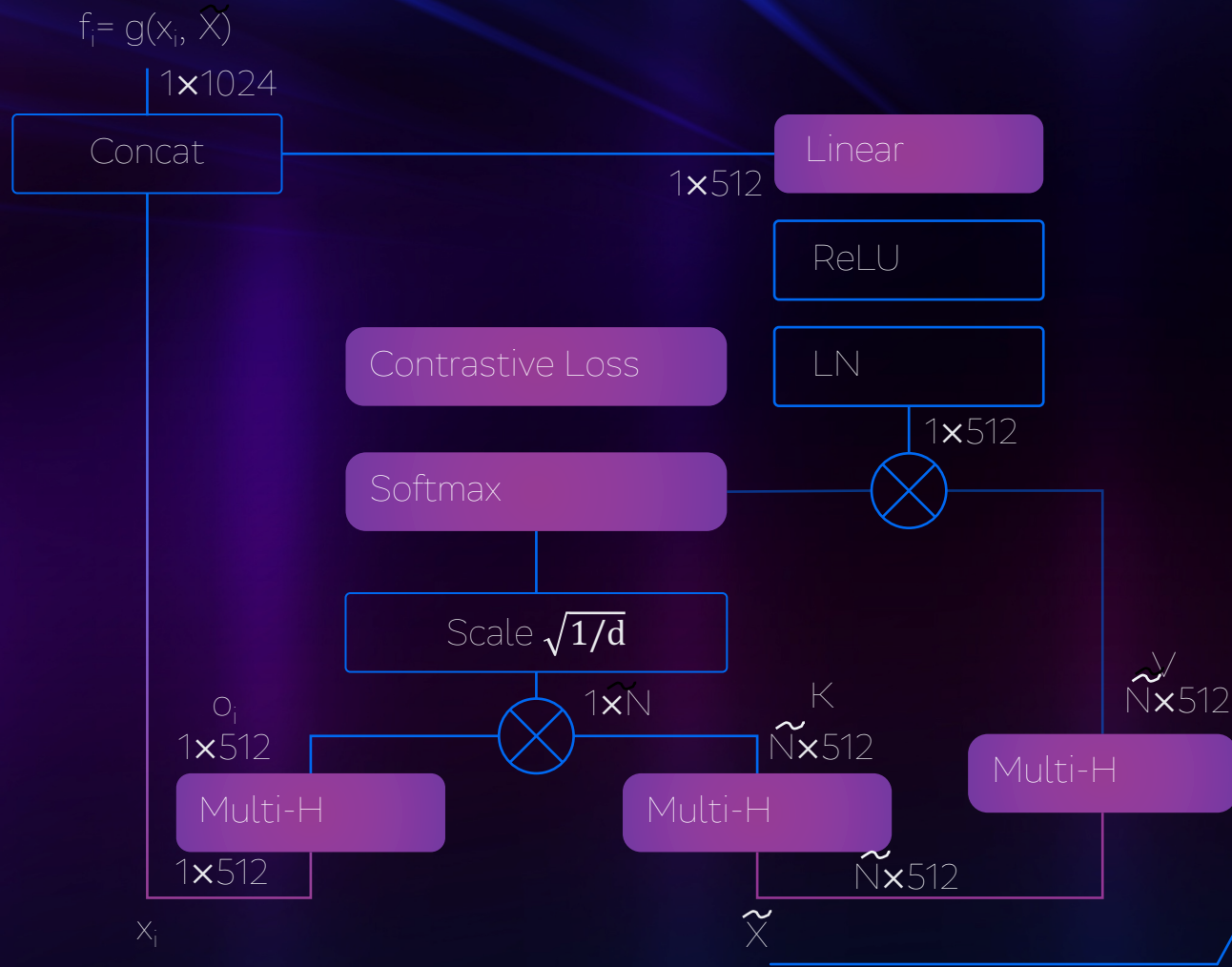


Фьюзинг модальностей

Bottleneck
transformer module



Contrastive attention



Postprocessing

