



# Уменьшаем задержку в live-трансляциях

Дмитрий Кравцов,  
Яндекс.Видеостриминг

Алексей Гусев,  
Яндекс.Дзен



# Где можно посмотреть видео через нас?

**КиноПоиск**

**ЯндексТВ**

Яндекс  Дзен

Яндекс  Алиса

Яндекс  Спорт

Яндекс  Музыка

Яндекс  Погода

Яндекс  Маркет

# Спортивные трансляции



# Стример и подписчики



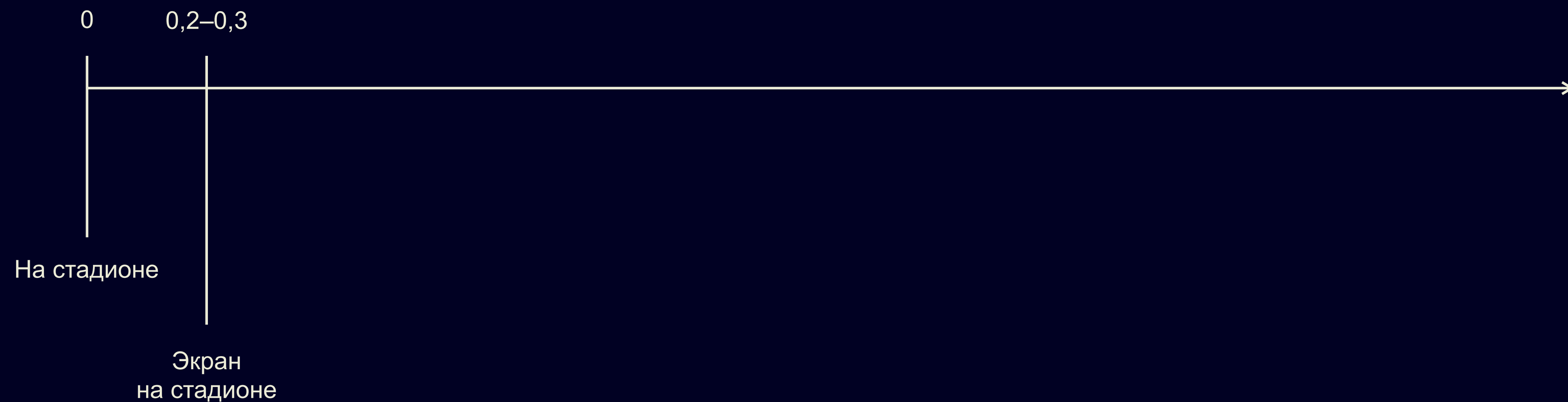
# Как можно смотреть футбольный матч?



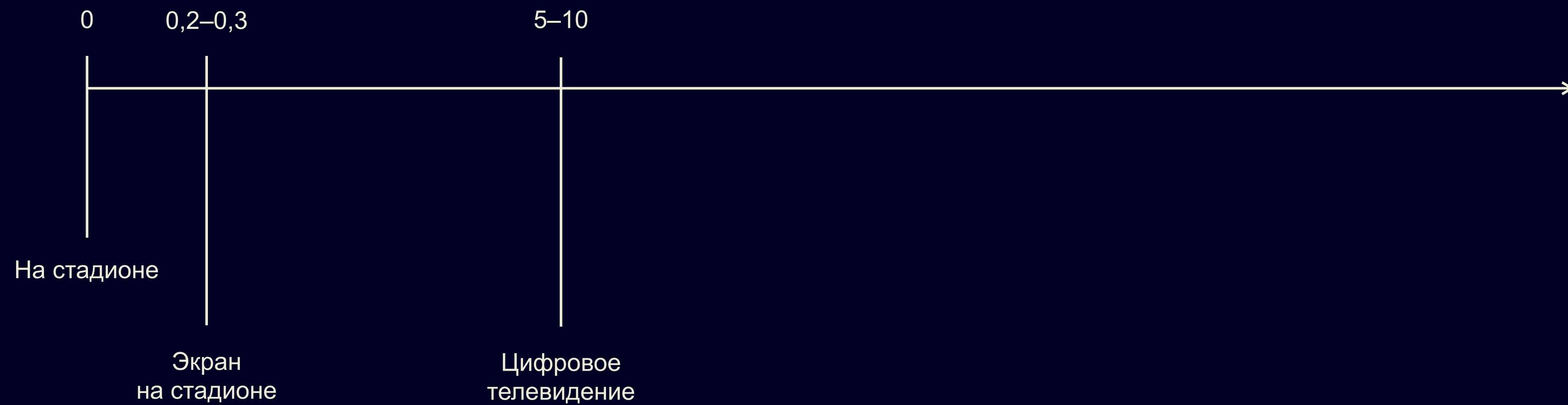
# Как можно смотреть футбольный матч?



# Как можно смотреть футбольный матч?

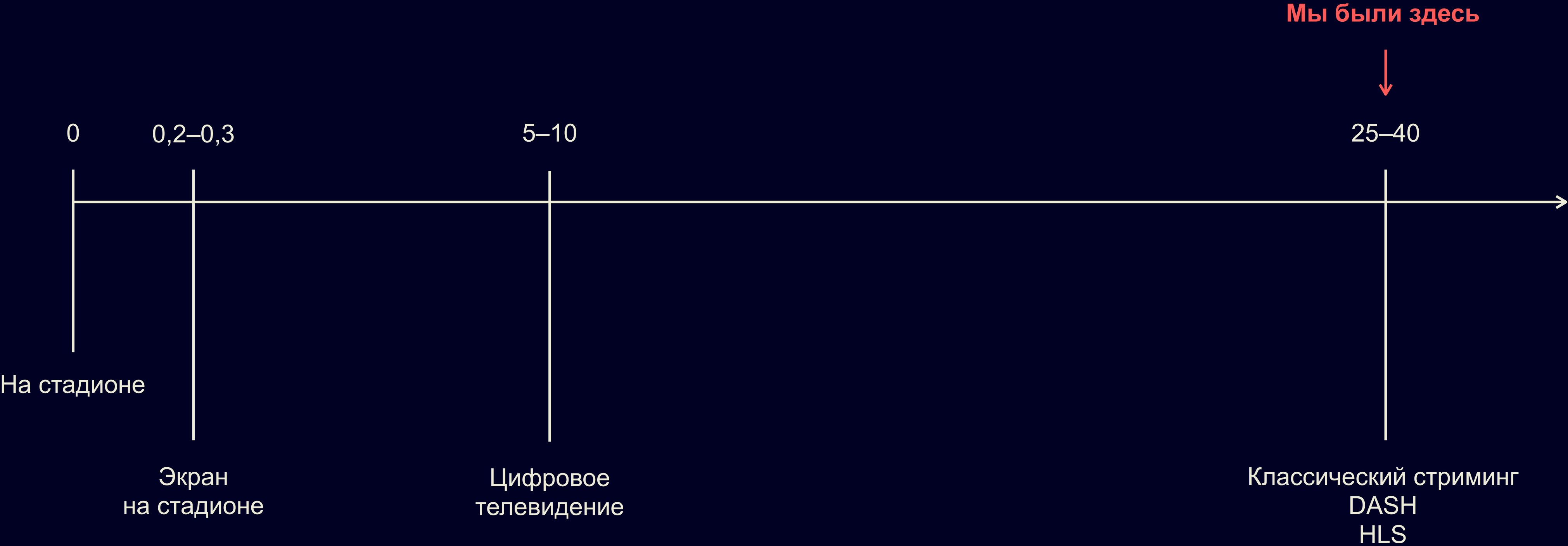


# Как можно смотреть футбольный матч?





# Как можно смотреть футбольный матч?



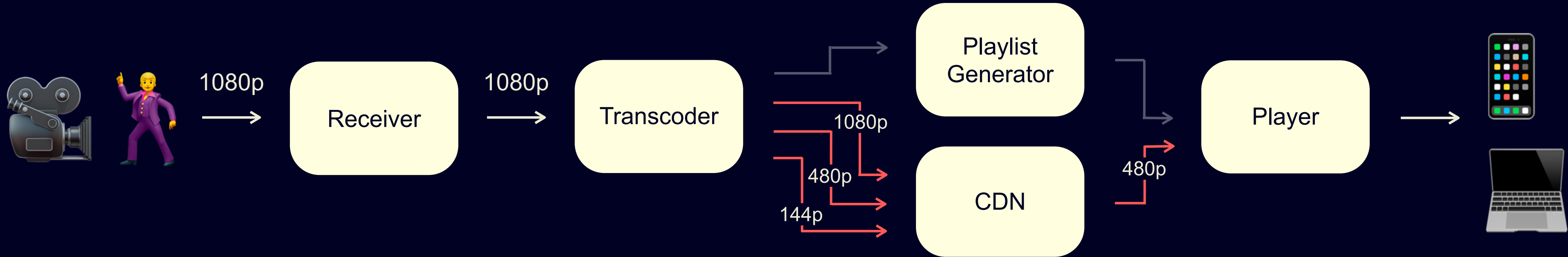
# Бюджет задержки

Во что мы целились?

10 секунд

# Старая архитектура

Как все работало раньше?



# Прием сигнала

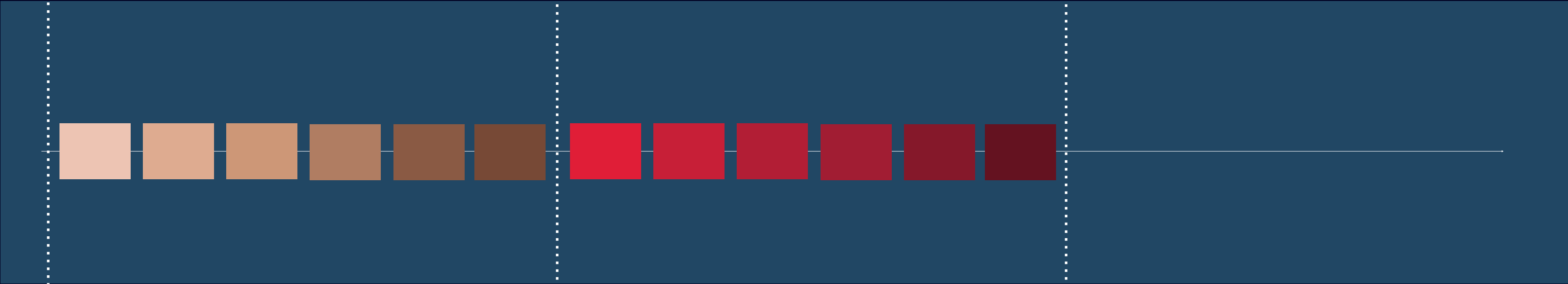


12:00:00

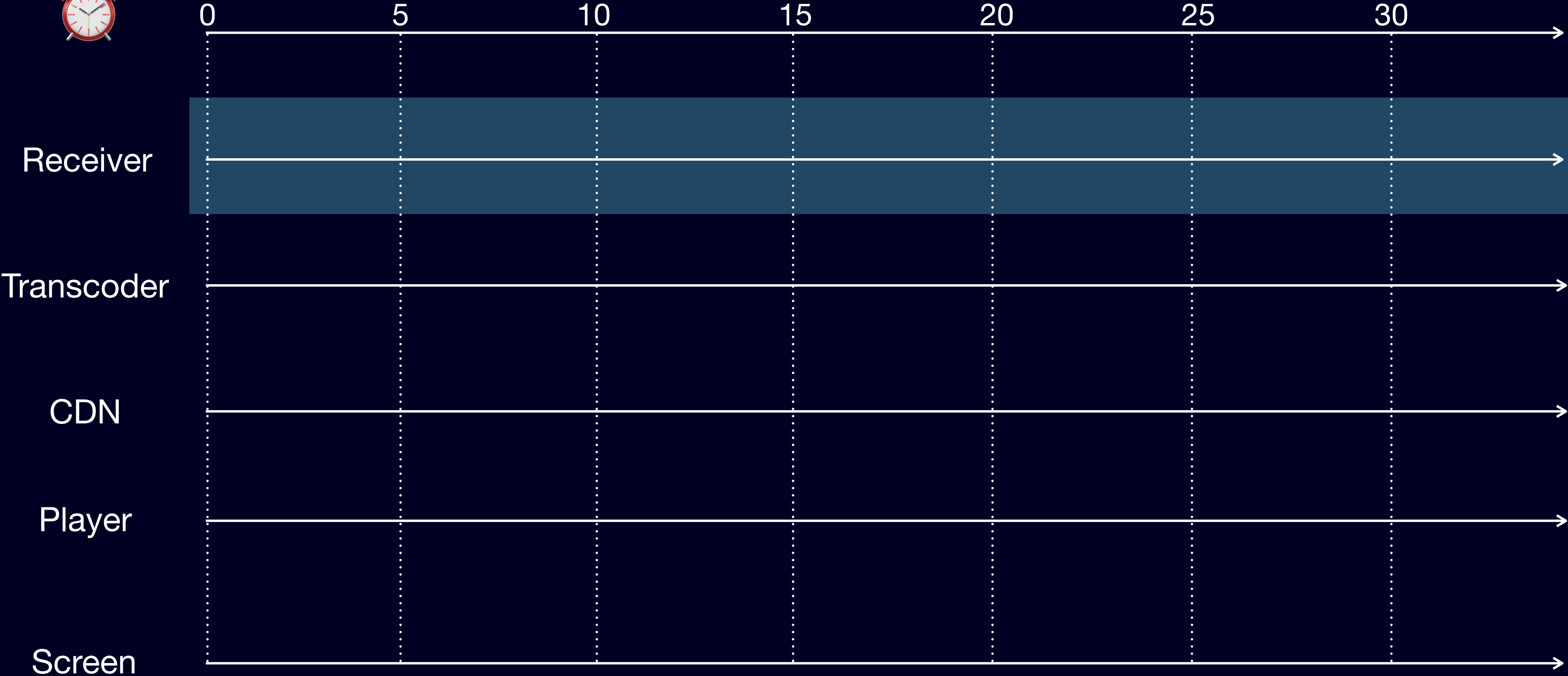
12:00:05

12:00:10

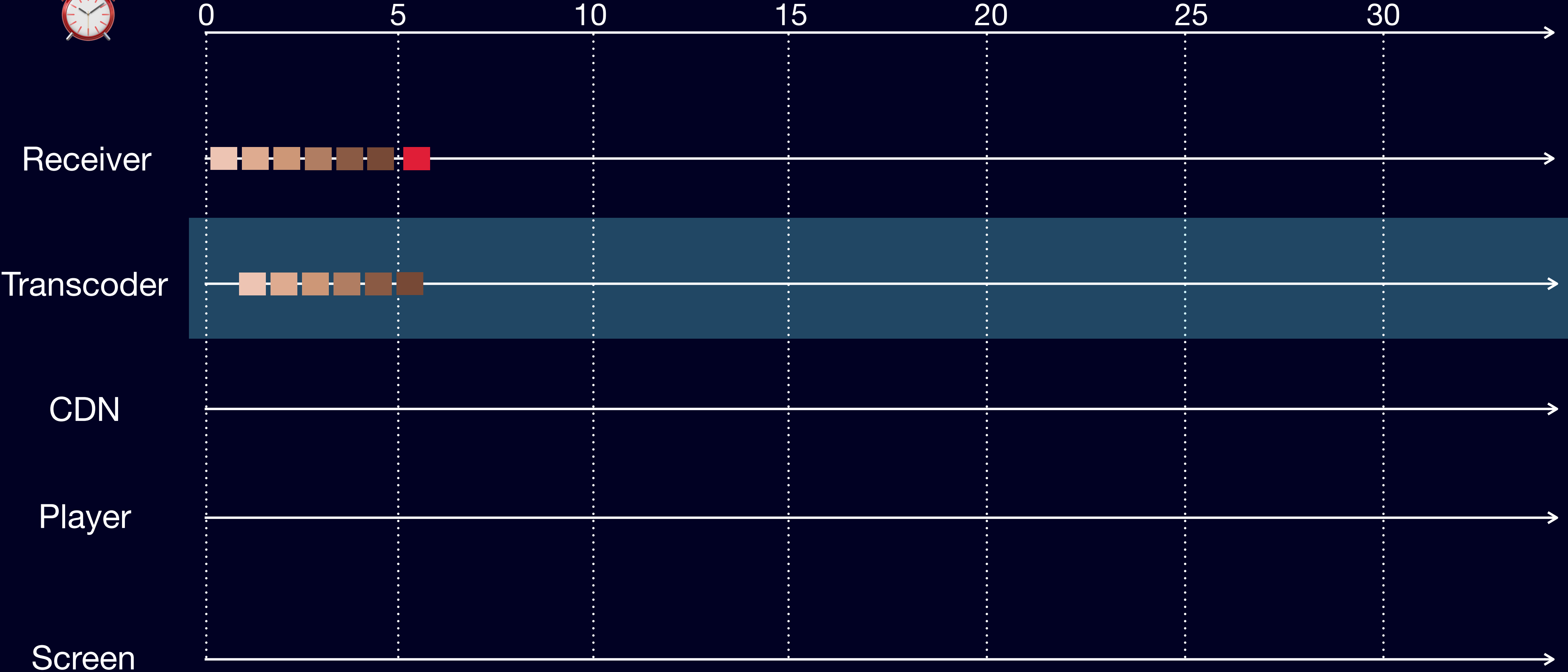
Receiver



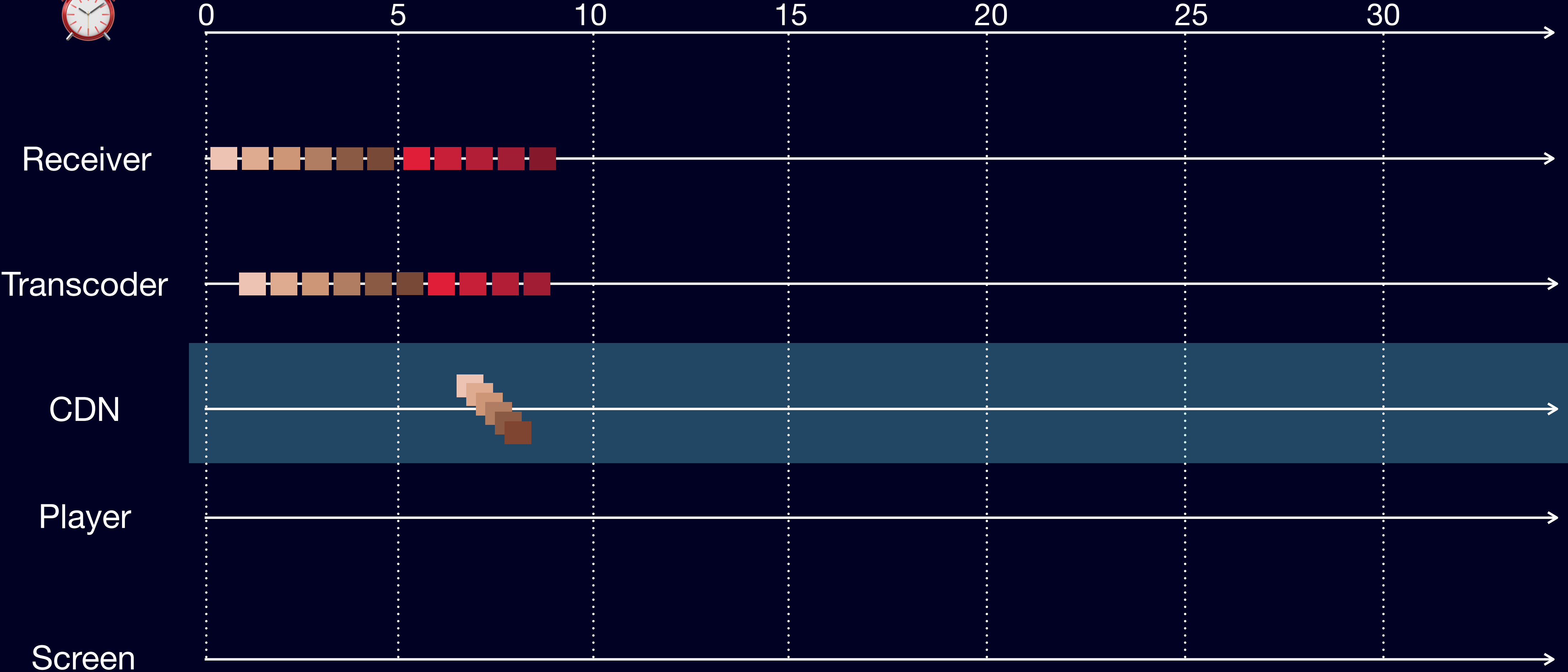
# Откуда берется задержка?



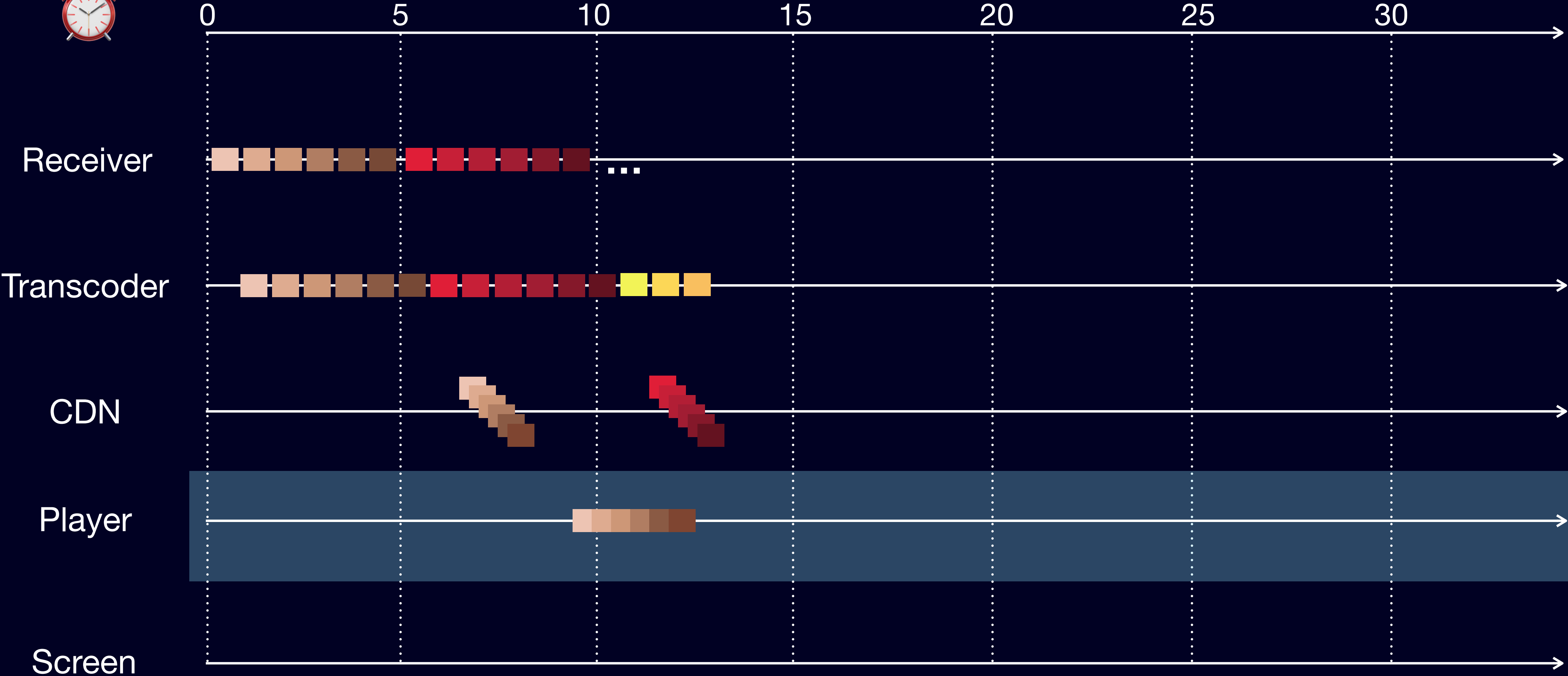
# Откуда берется задержка?



# Откуда берется задержка?

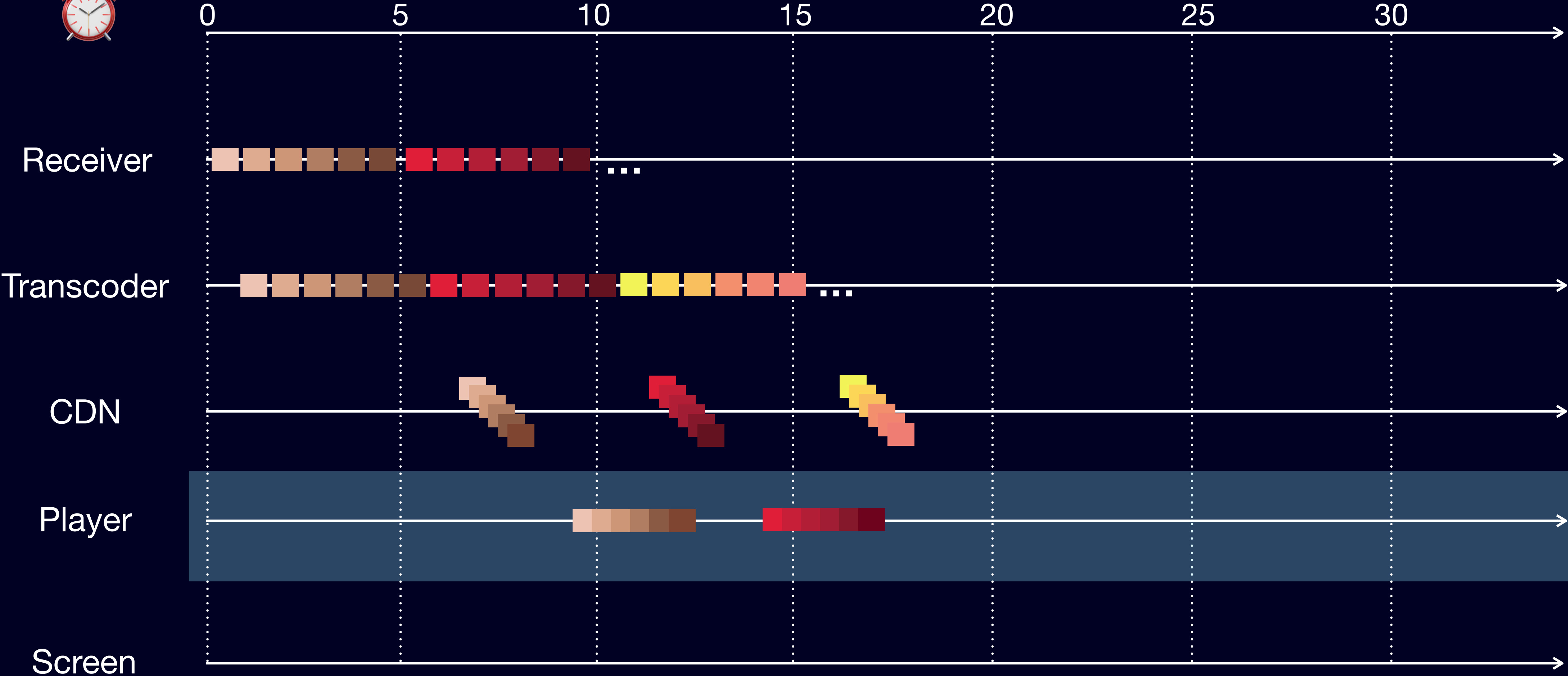


# Откуда берется задержка?

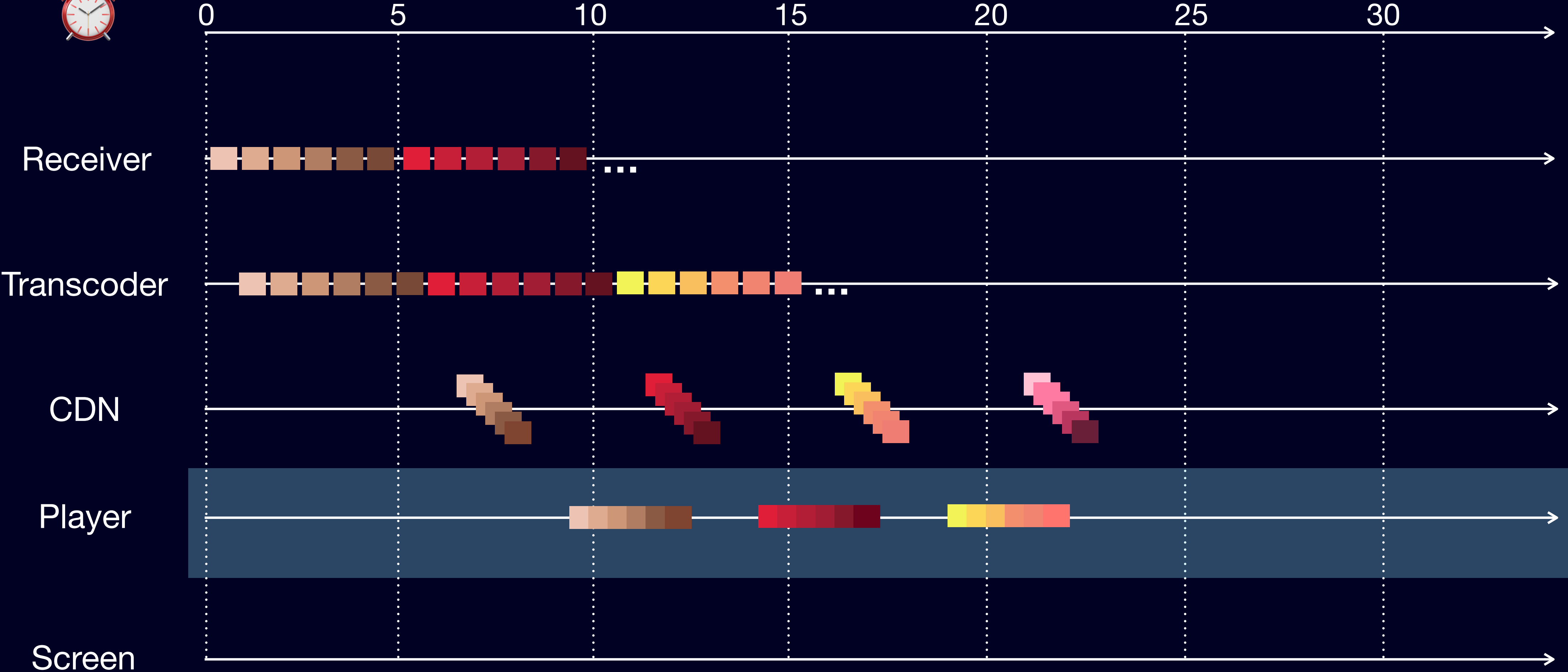




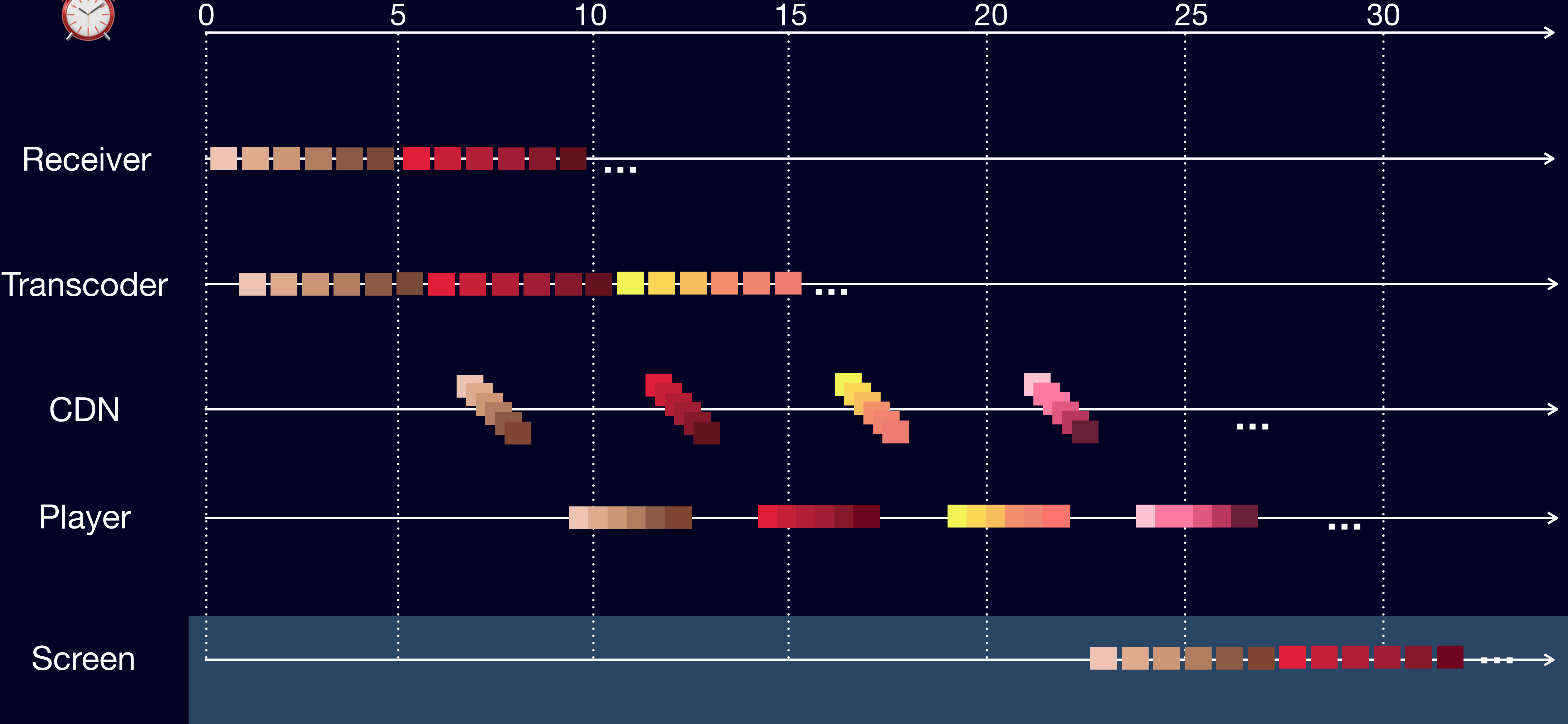
# Откуда берется задержка?



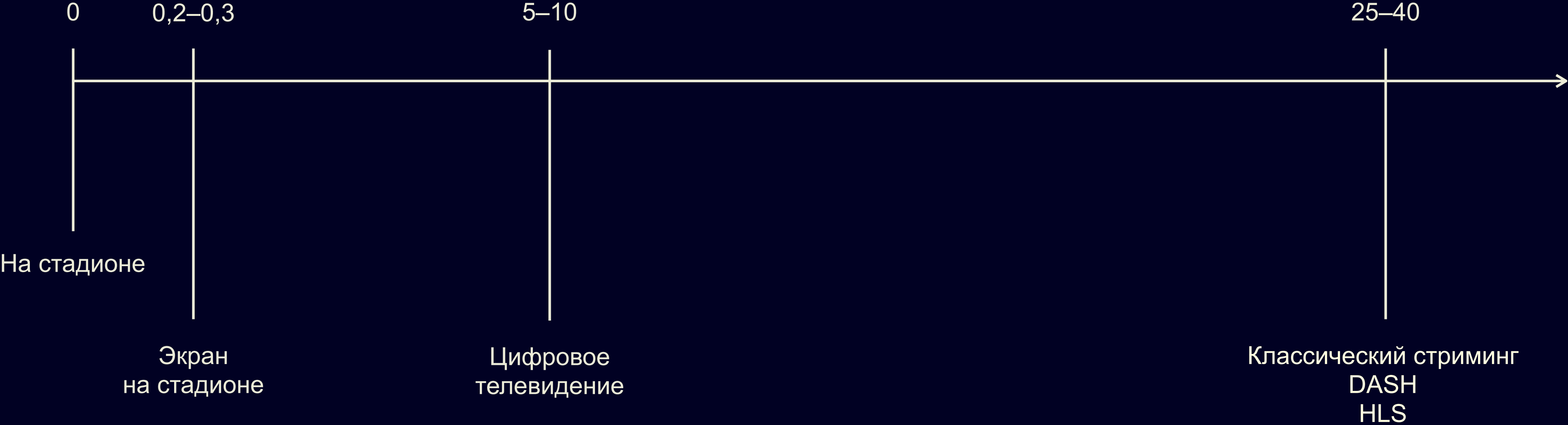
# Откуда берется задержка?



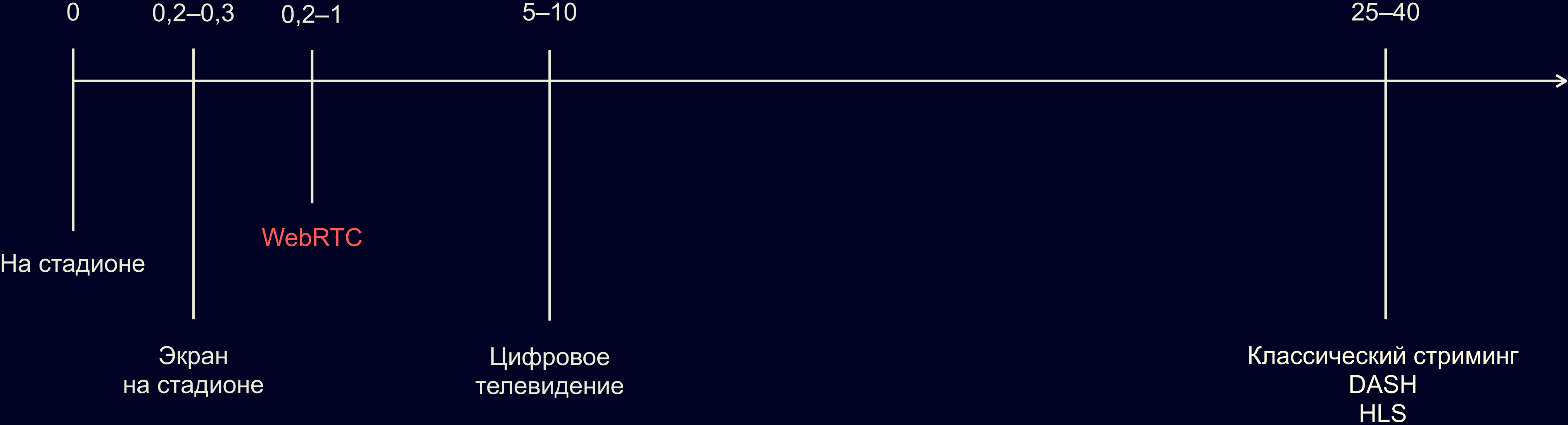
# Откуда берется задержка?



# Куда можно было двигаться



# Куда можно было двигаться



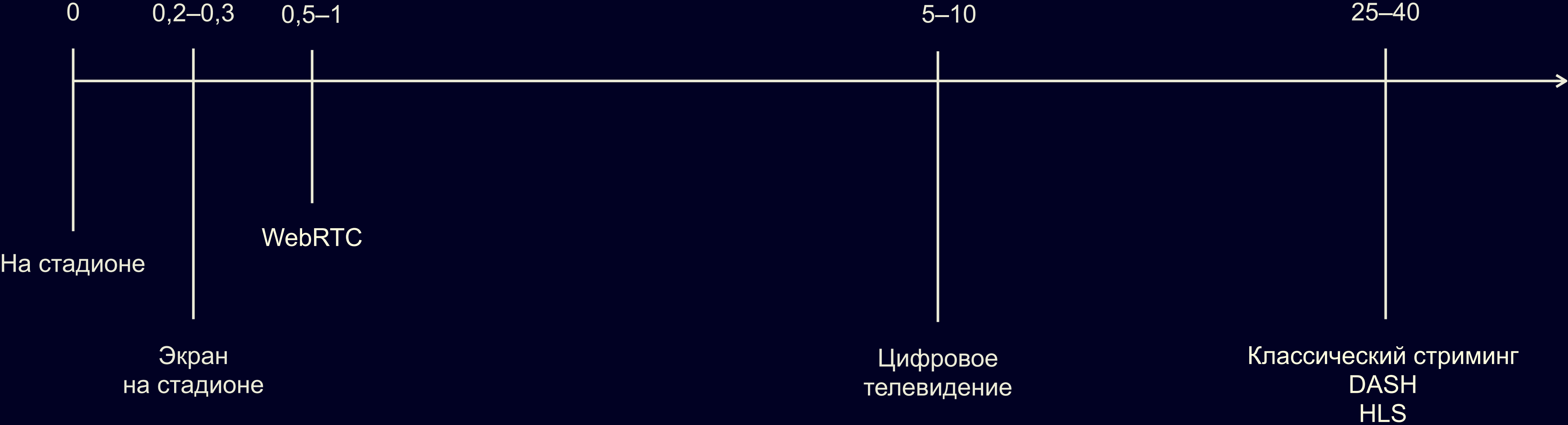
# WebRTC

+ Супер-маленькая  
задержка

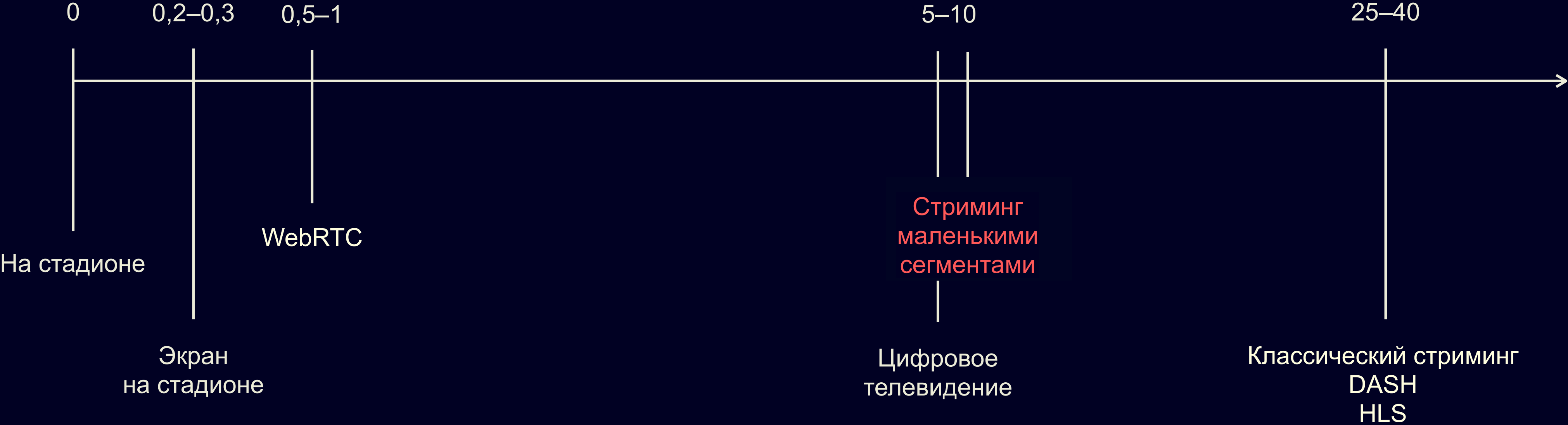
- Отсутствует гарантия  
на доставку

- Не все клиенты  
поддерживают

# Куда можно было двигаться

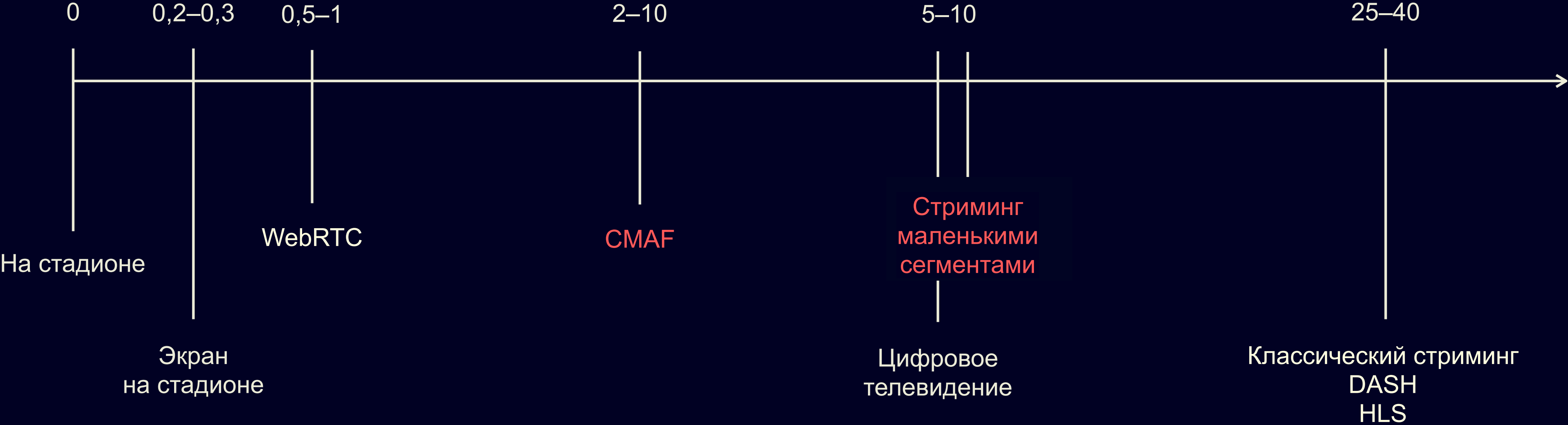


# Куда можно было двигаться

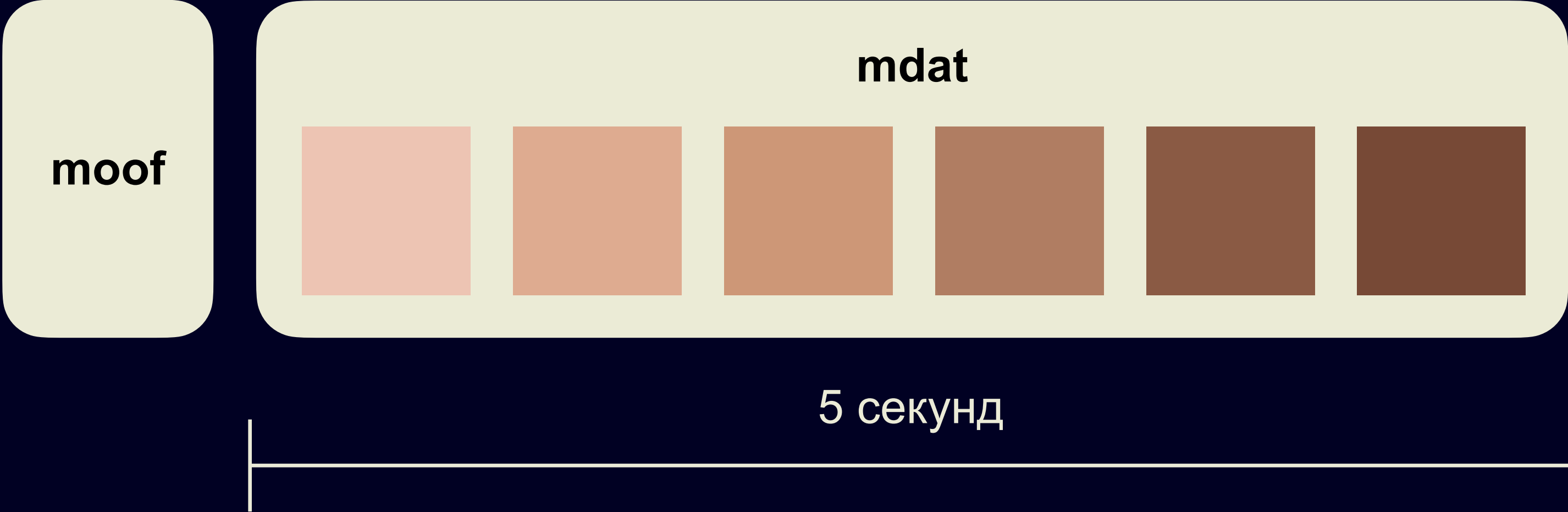




# Куда можно было двигаться

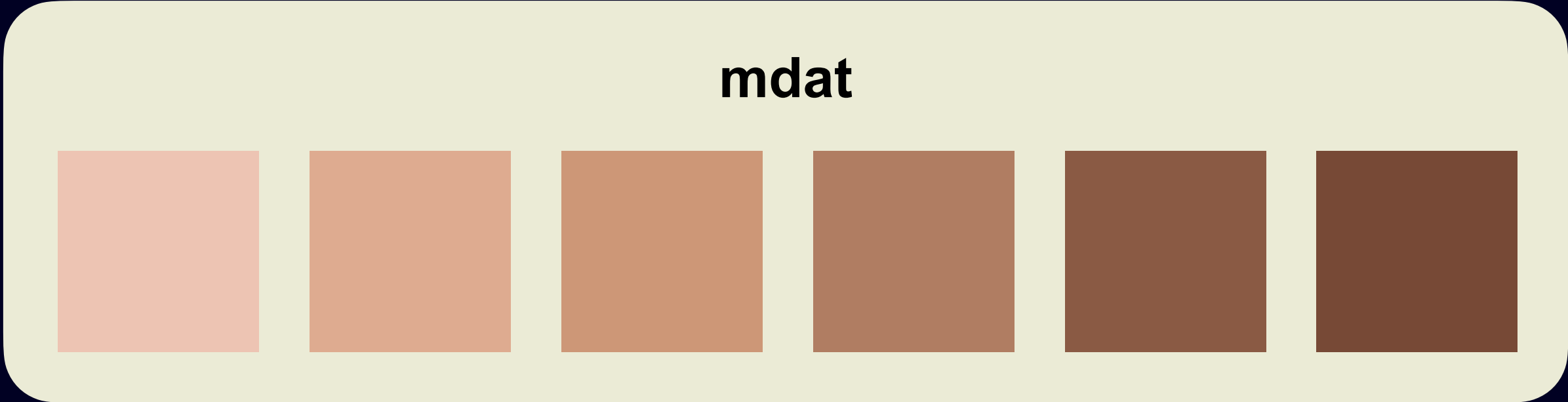


# Как устроен сегмент?

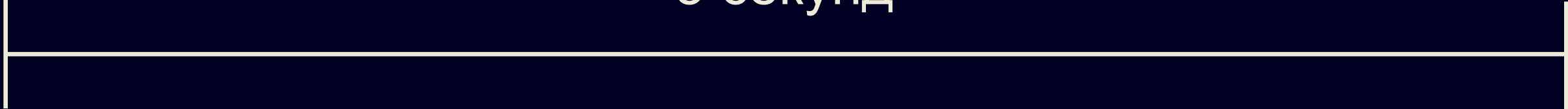


# Как устроен сегмент?

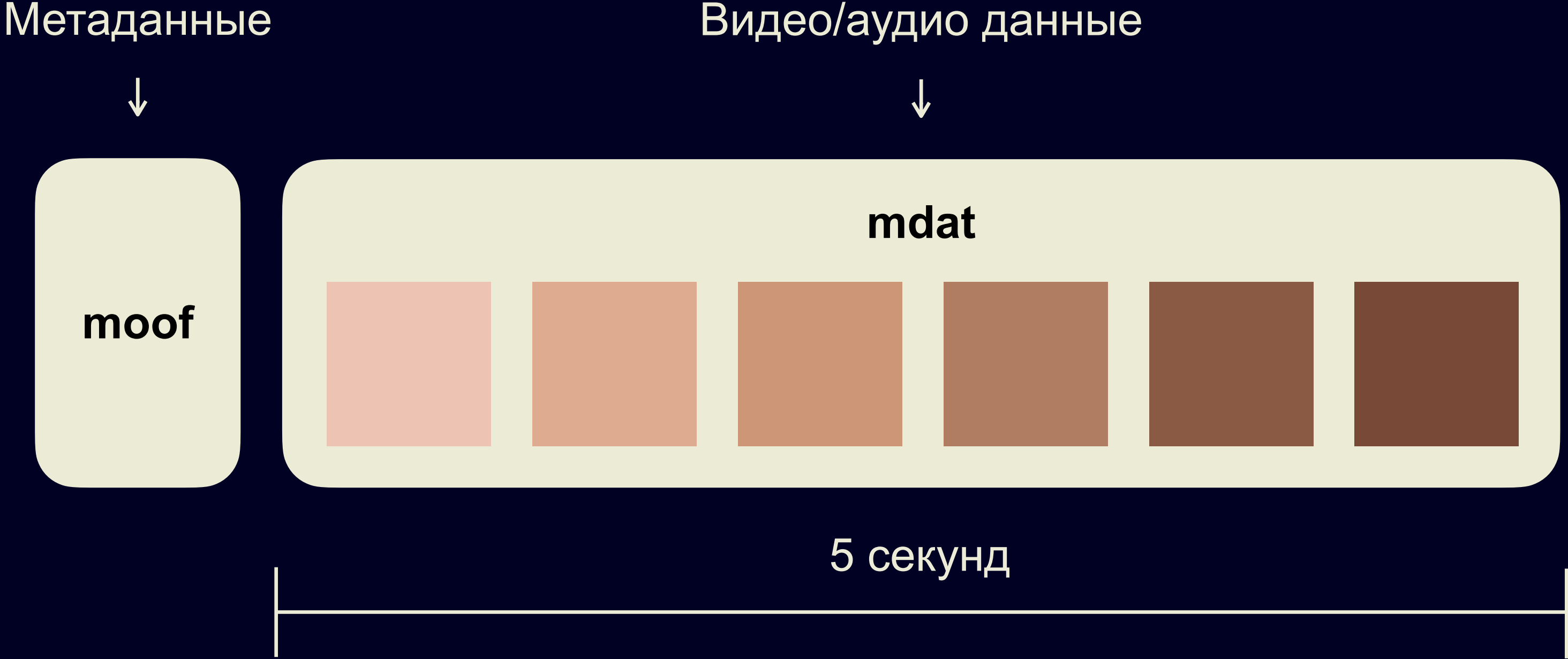
Метаданные



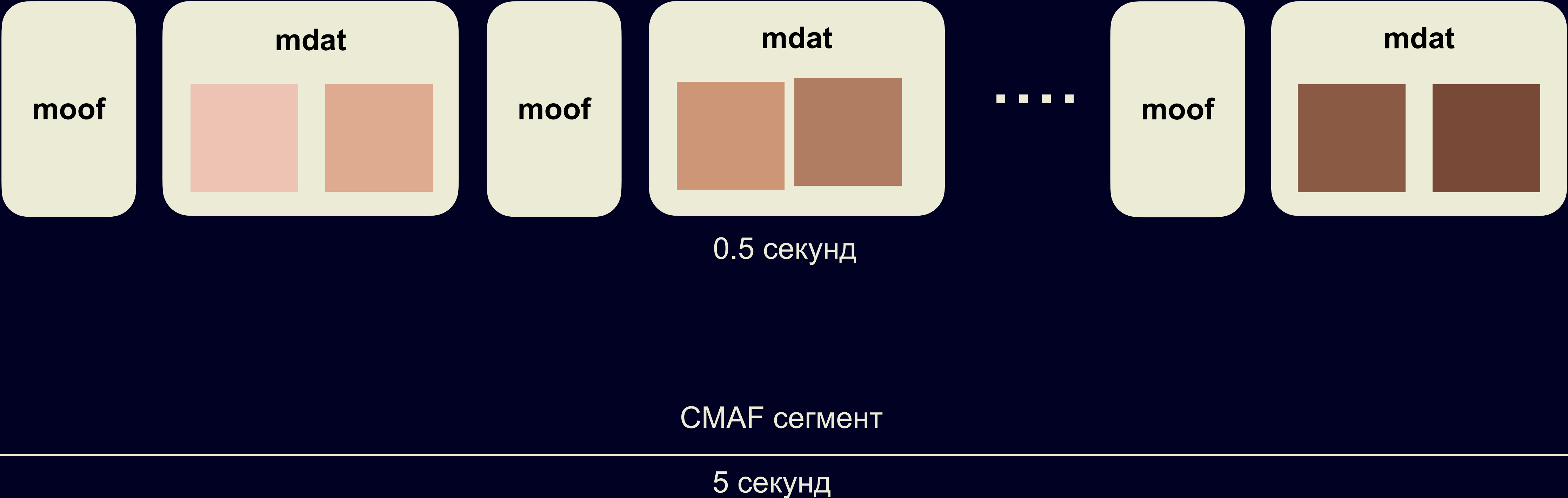
5 секунд



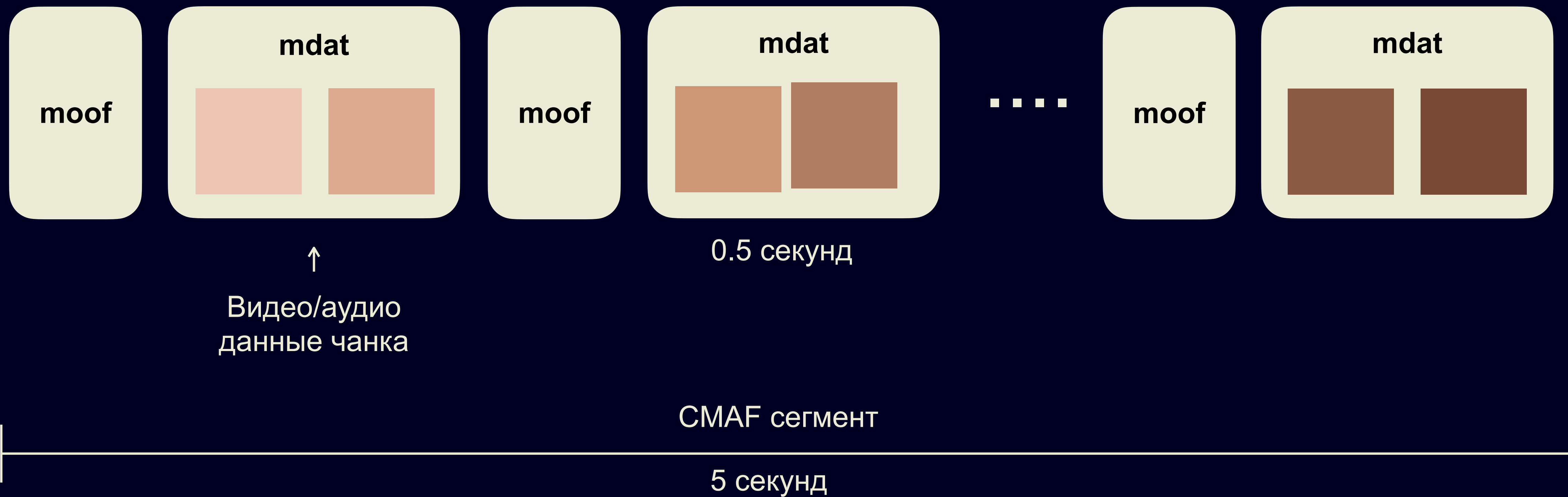
# Как устроен сегмент?



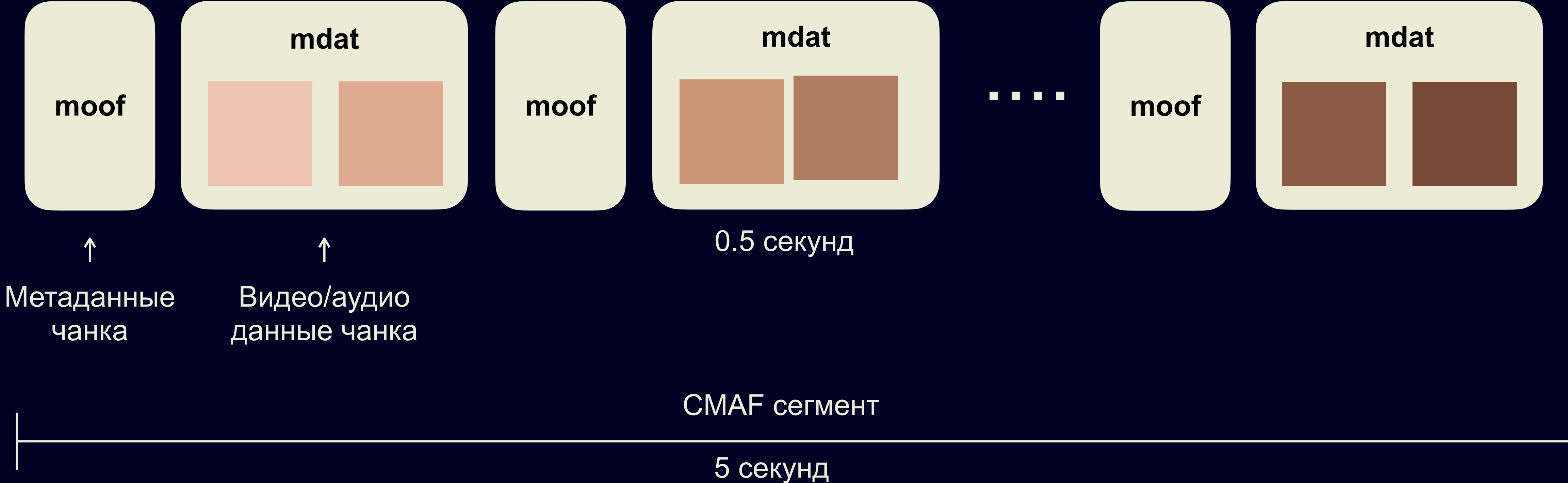
# Как устроен CMAF-сегмент?



# Как устроен CMAF-сегмент?



# Как устроен CMAF-сегмент?



# Как раздать неполный сегмент?

**Сегмент**



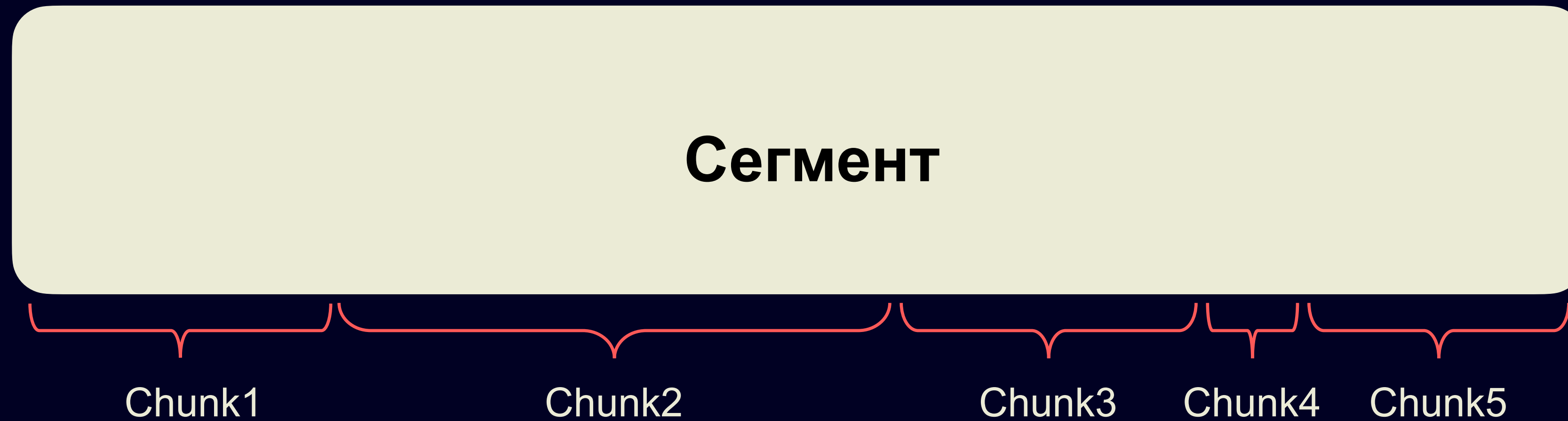
# Как раздать неполный сегмент?

Transfer-Encoding: chunked

**Сегмент**

# Как раздать неполный сегмент?

Transfer-Encoding: chunked



# Как раздать неполный сегмент?

Transfer-Encoding: chunked



Chunk1

Chunk2

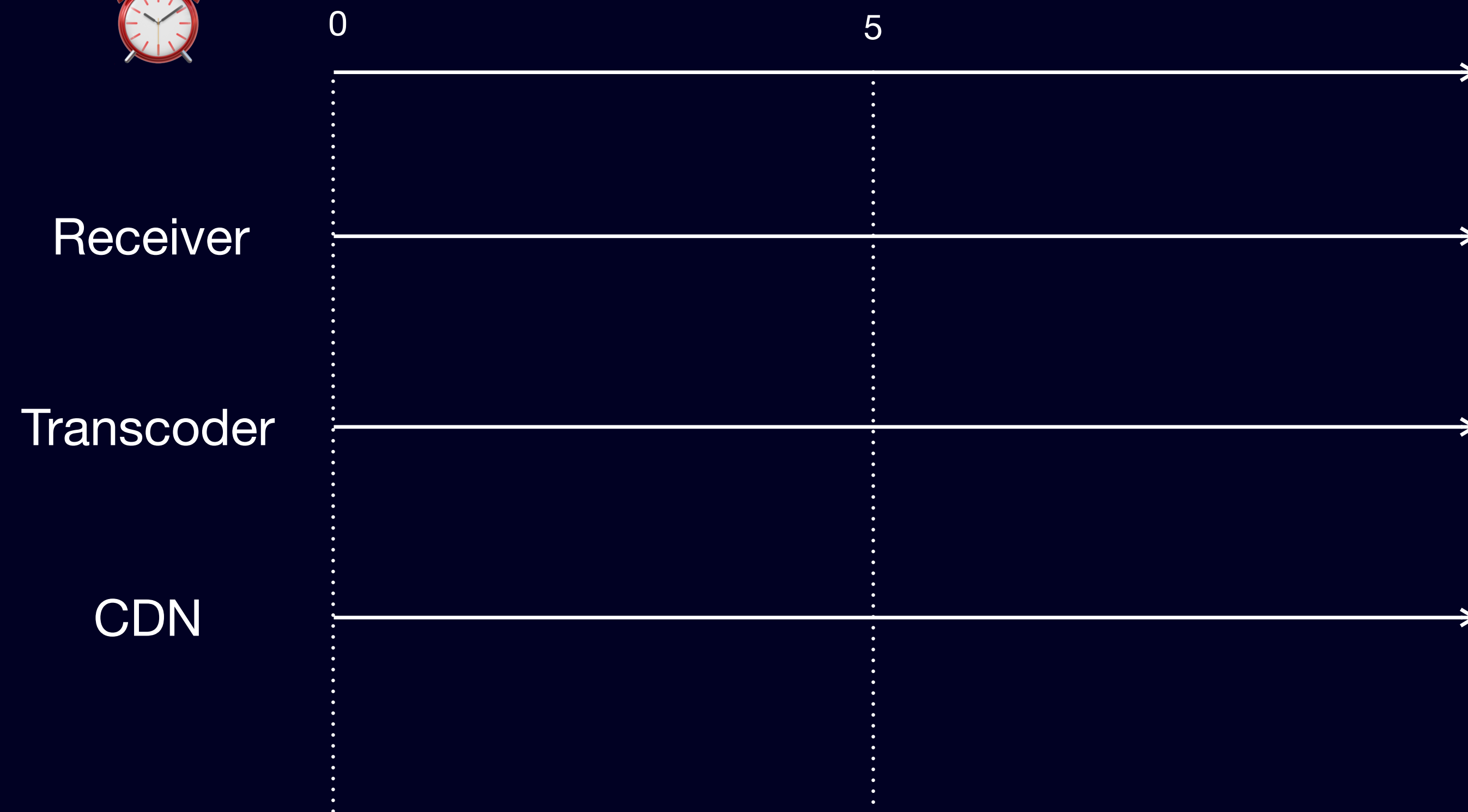
Chunk3

Chunk4

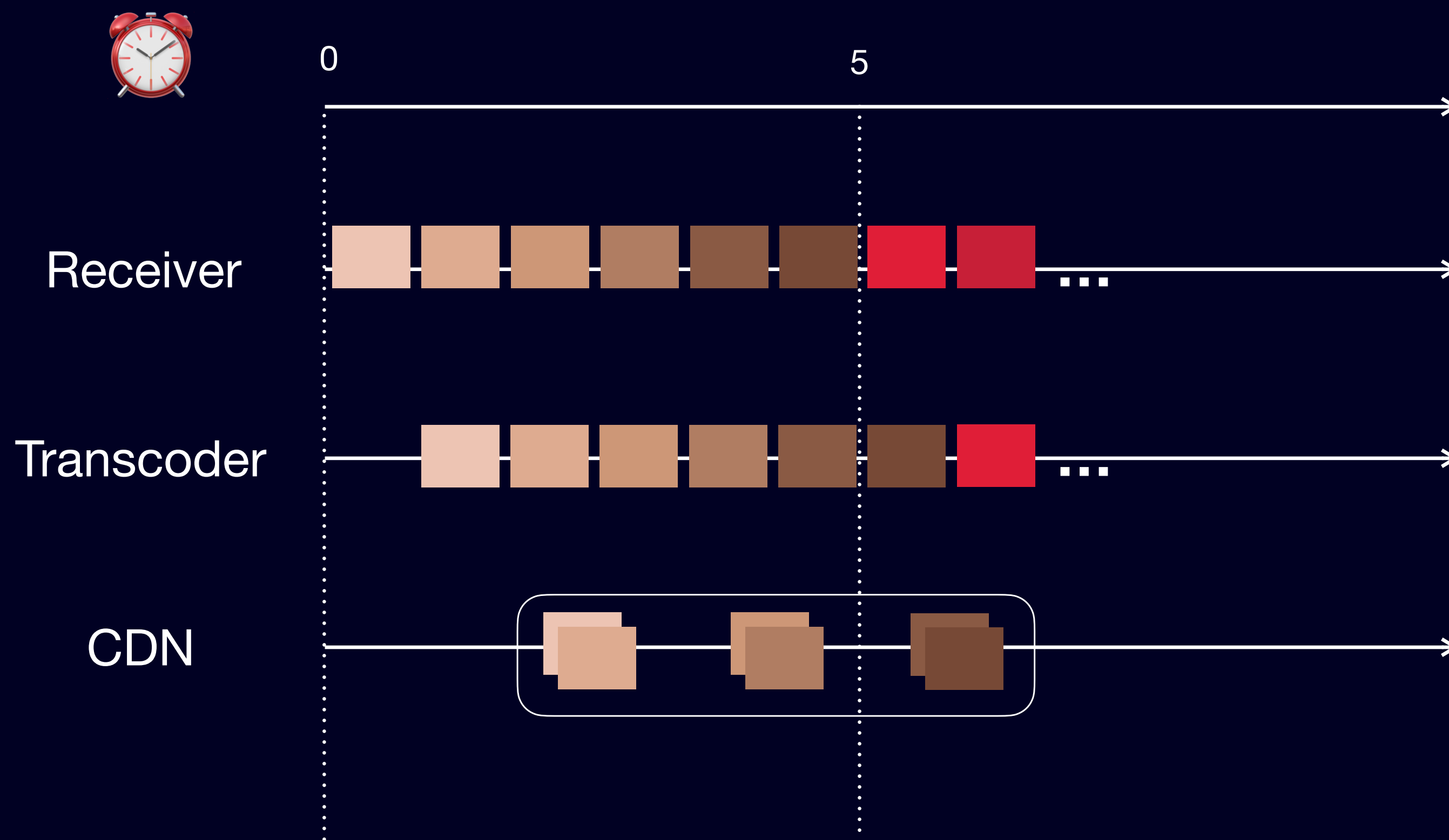
Chunk5

Empty chunk

# Задержка для СМАФ

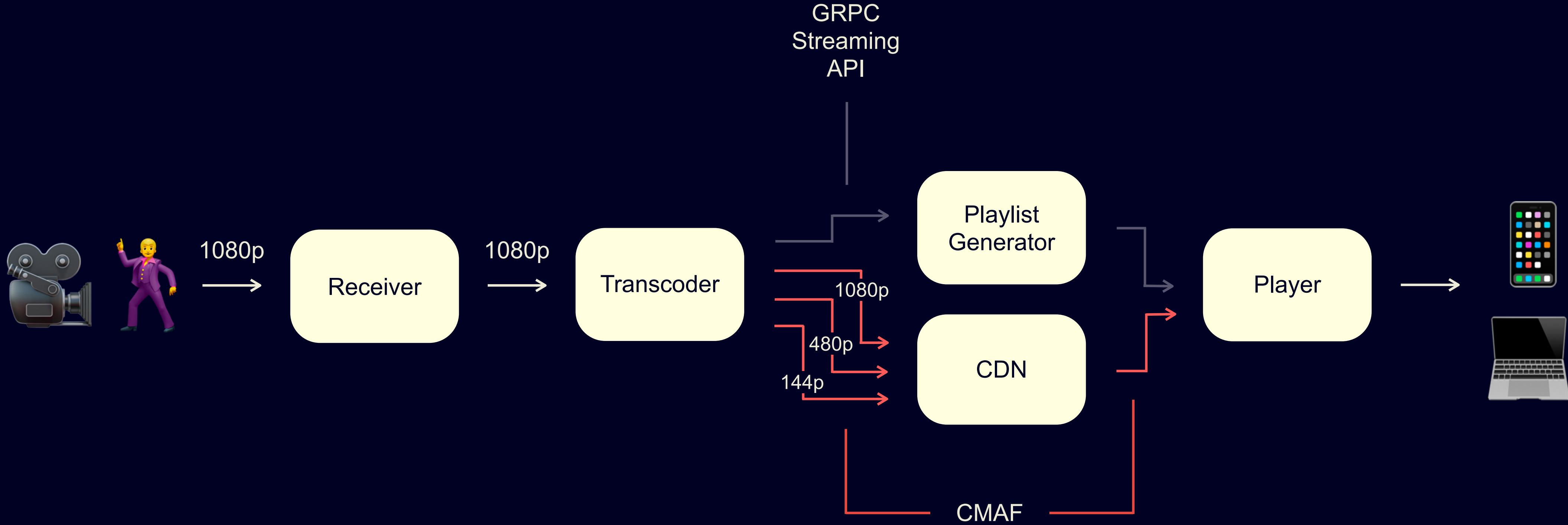


# Задержка для СМАФ

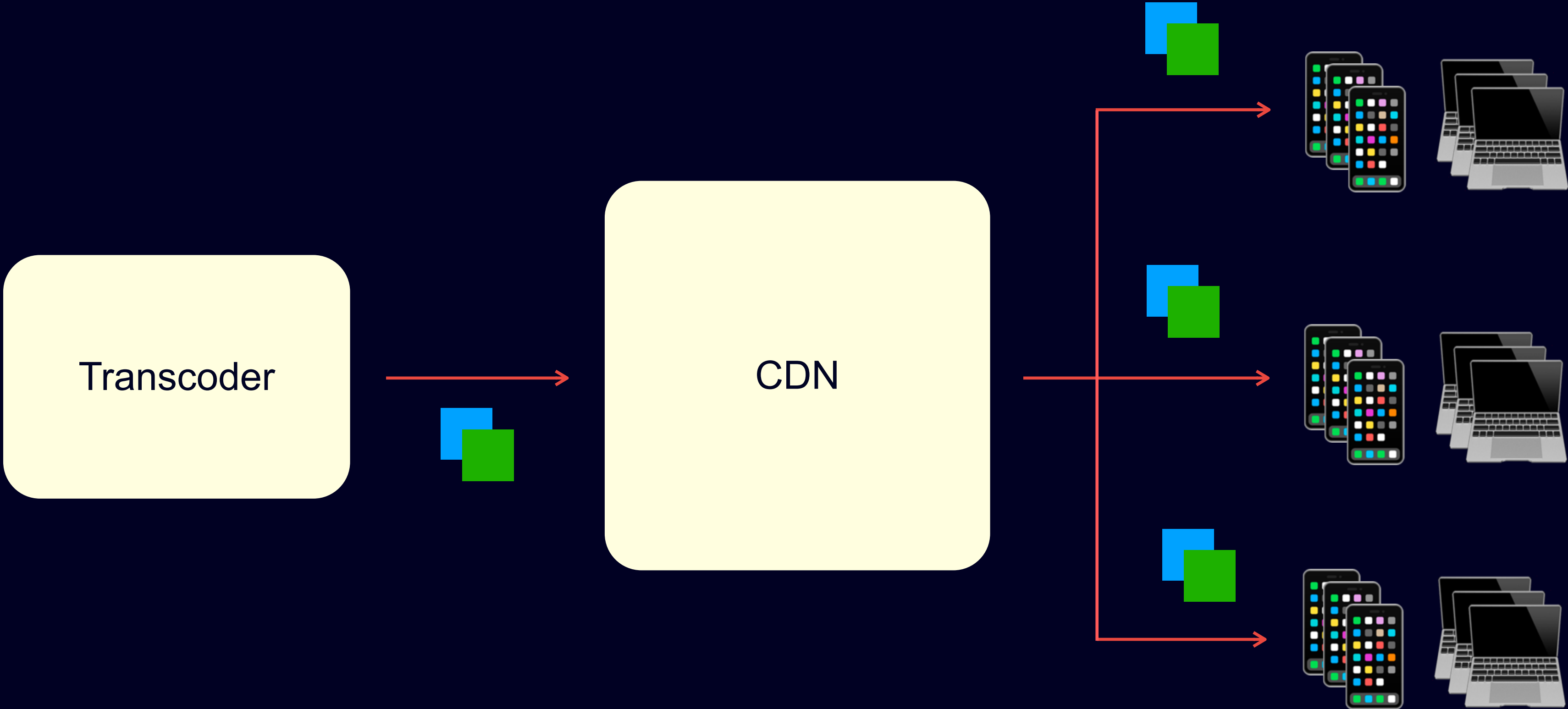


# Новая архитектура

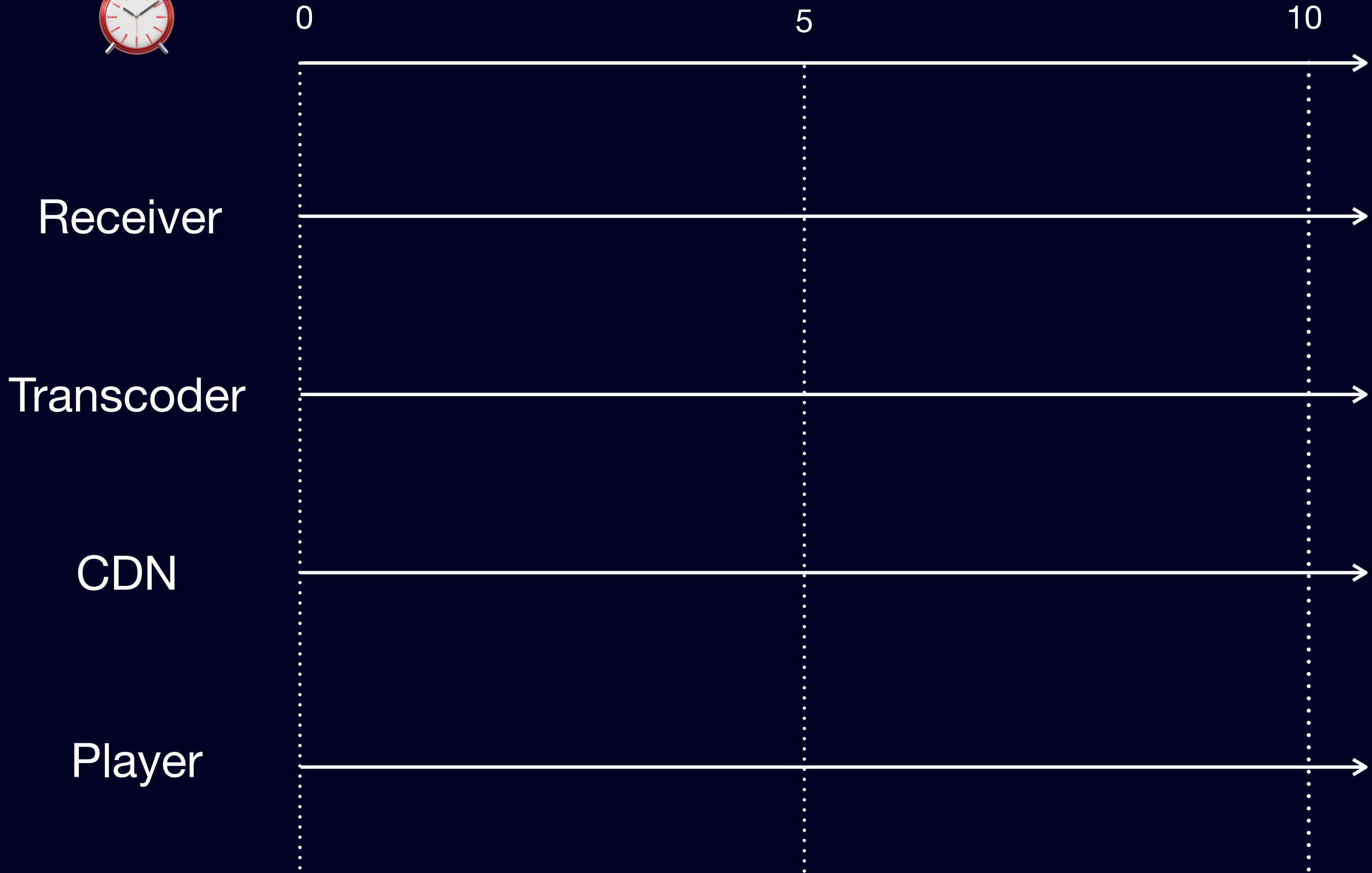
А что поменялось?



# Cache follow B CDN

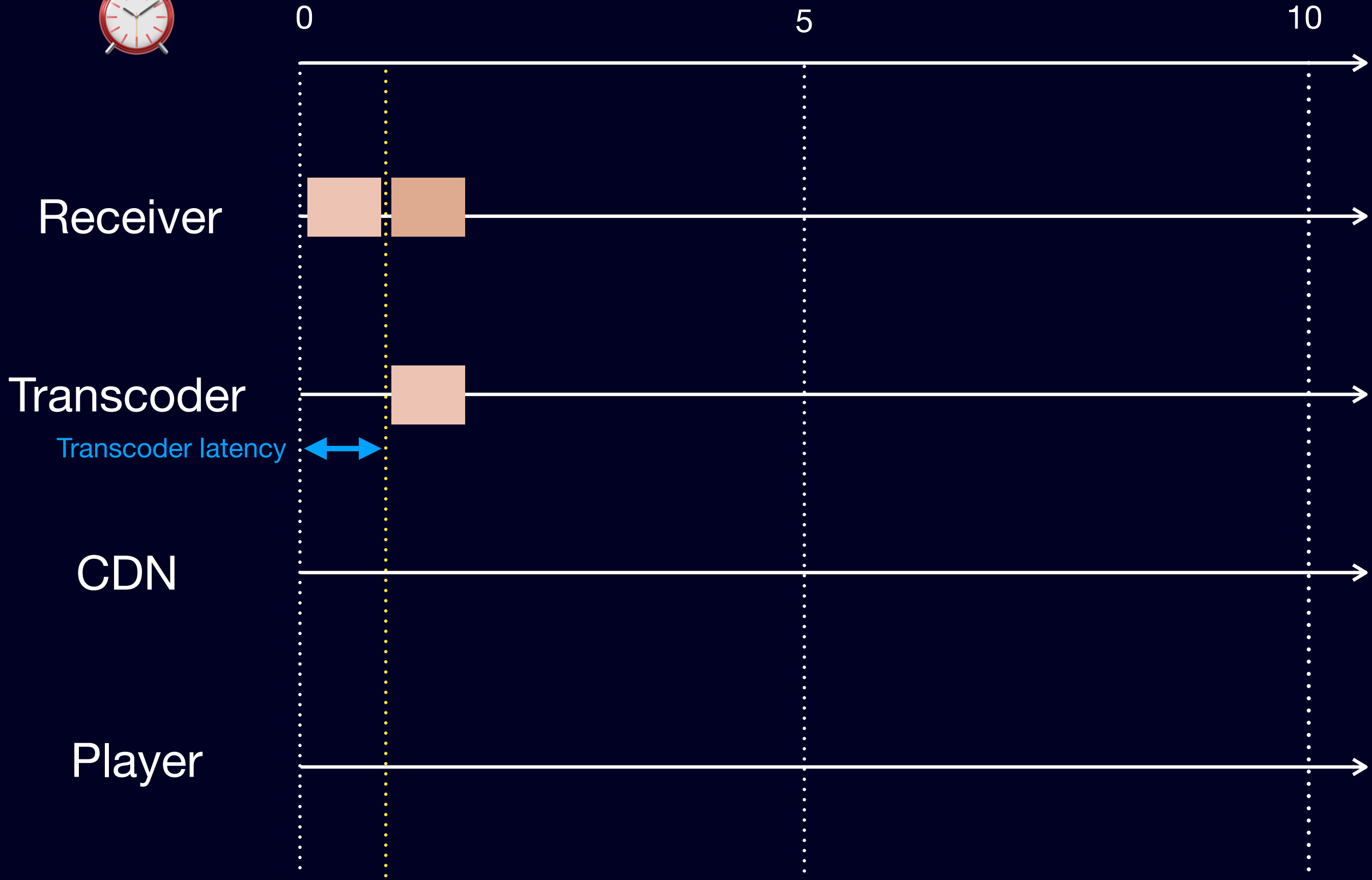


# Метрики

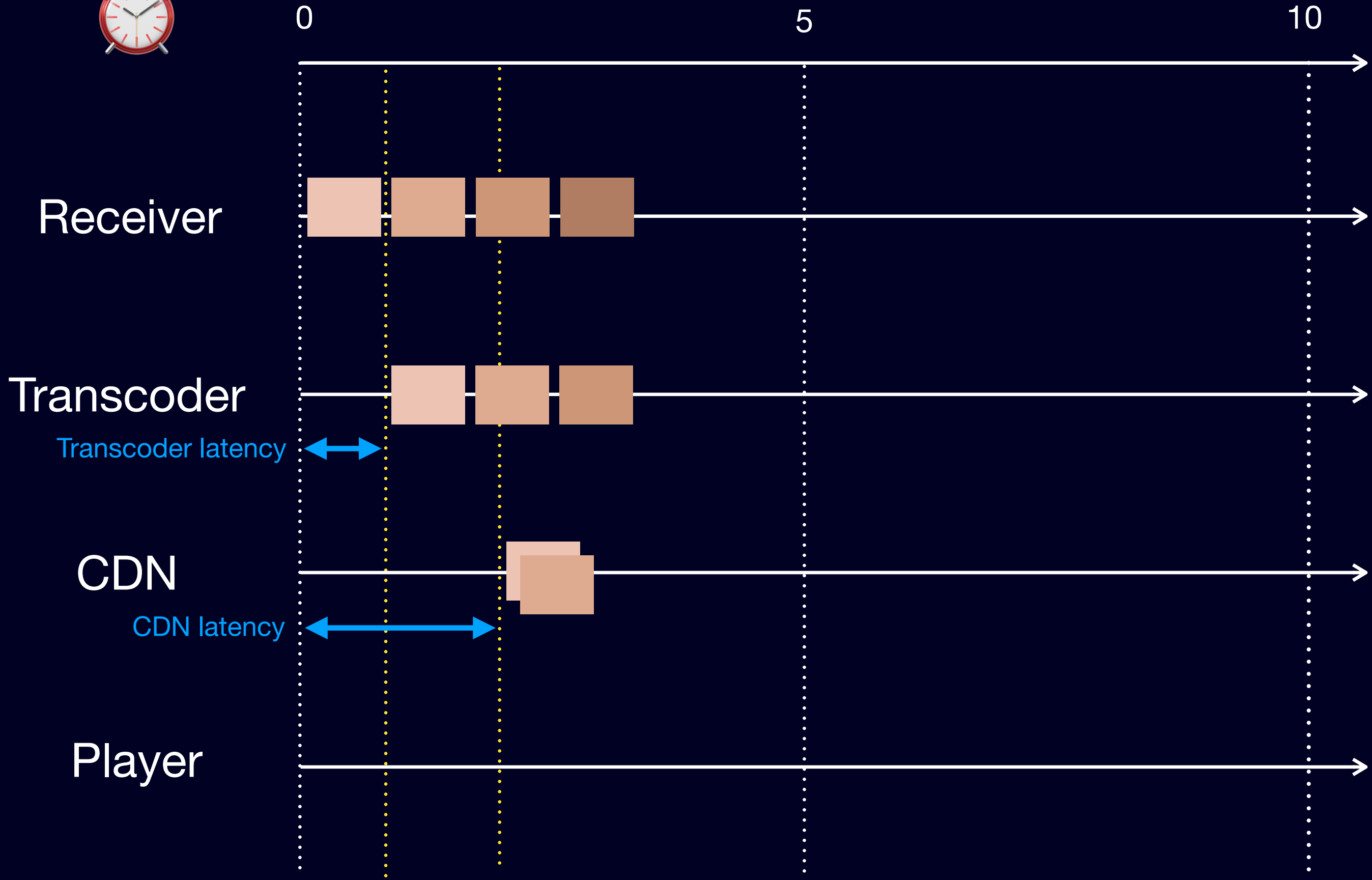




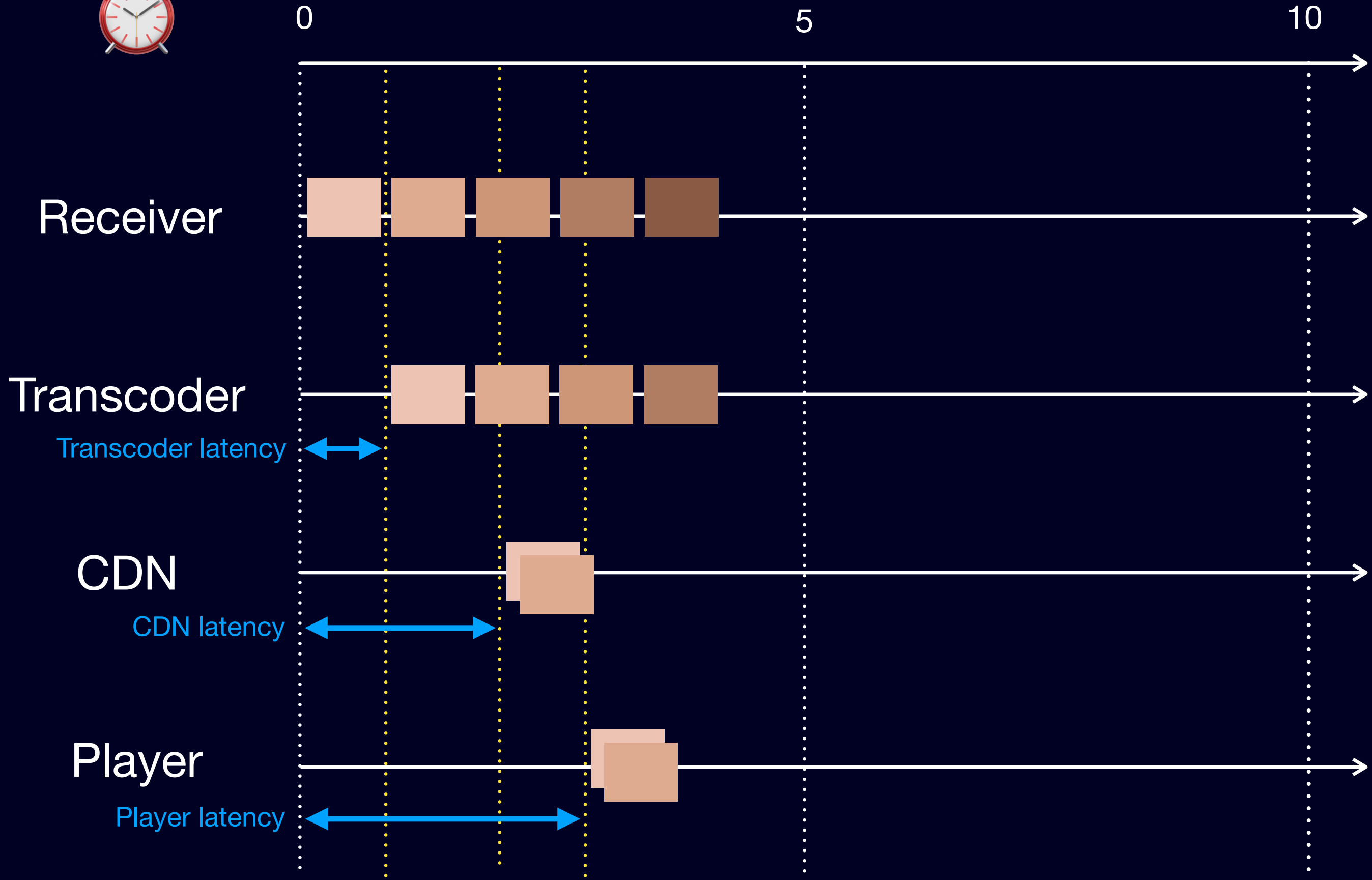
# Метрики



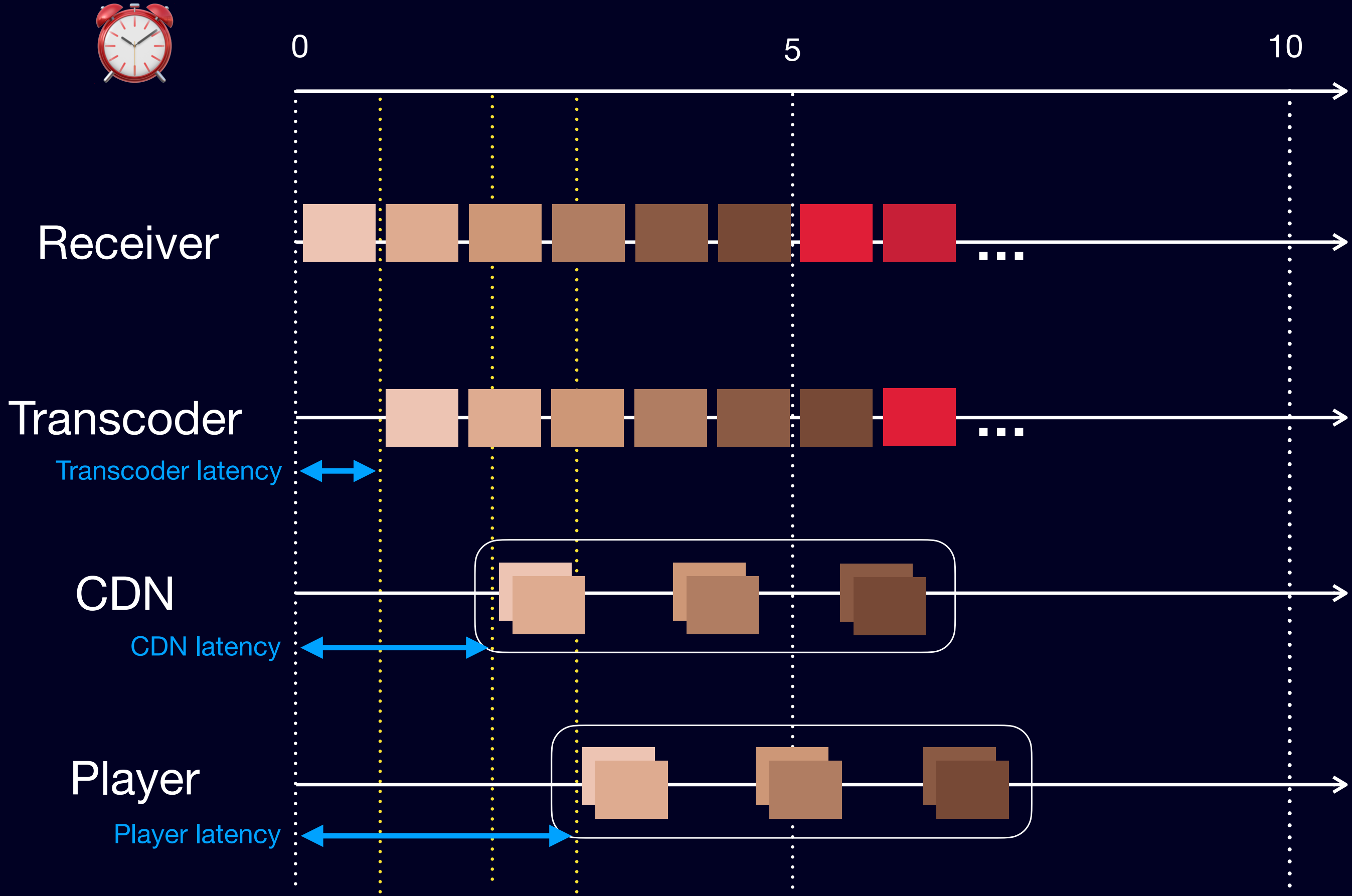
# Метрики



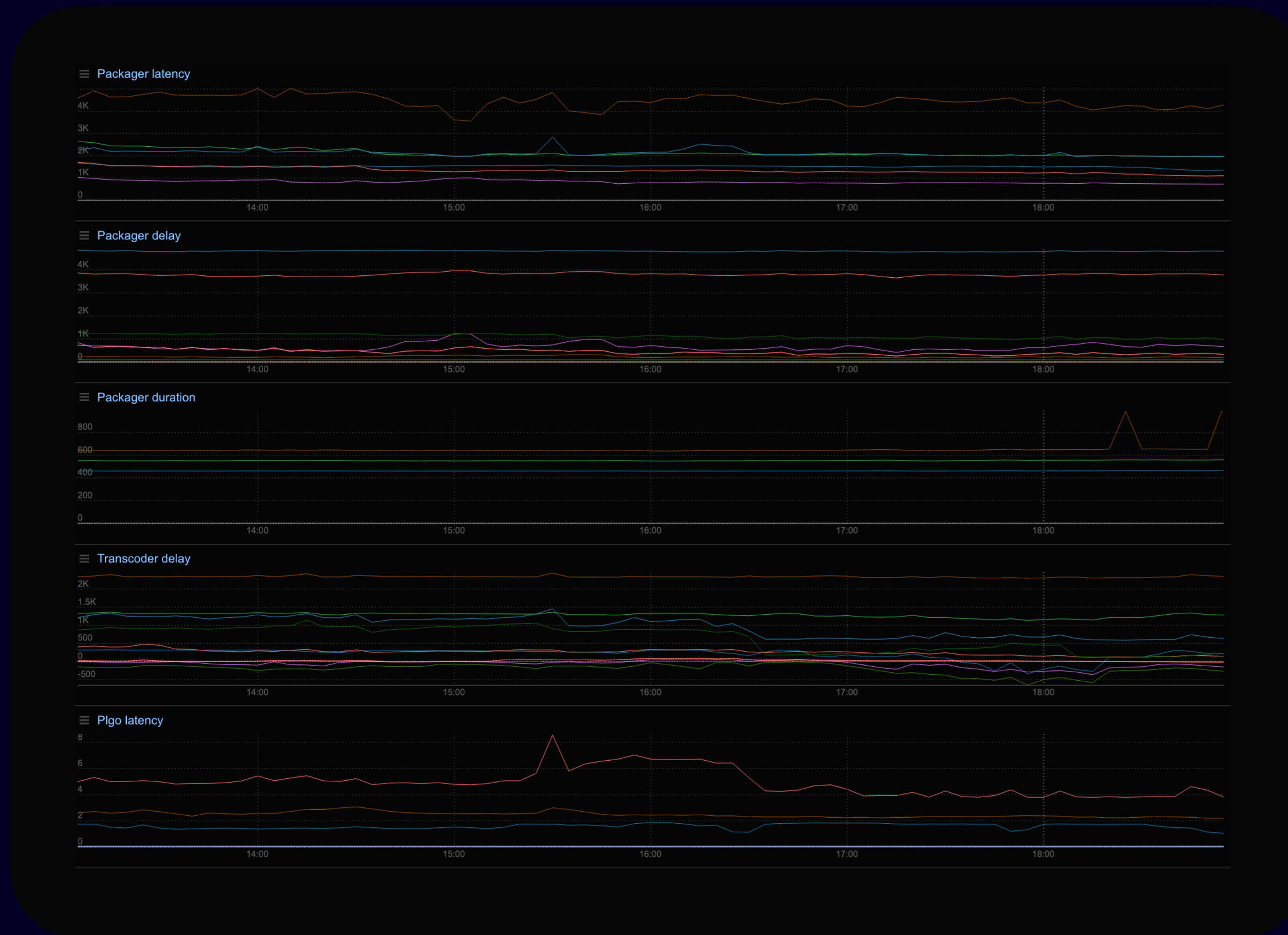
# Метрики



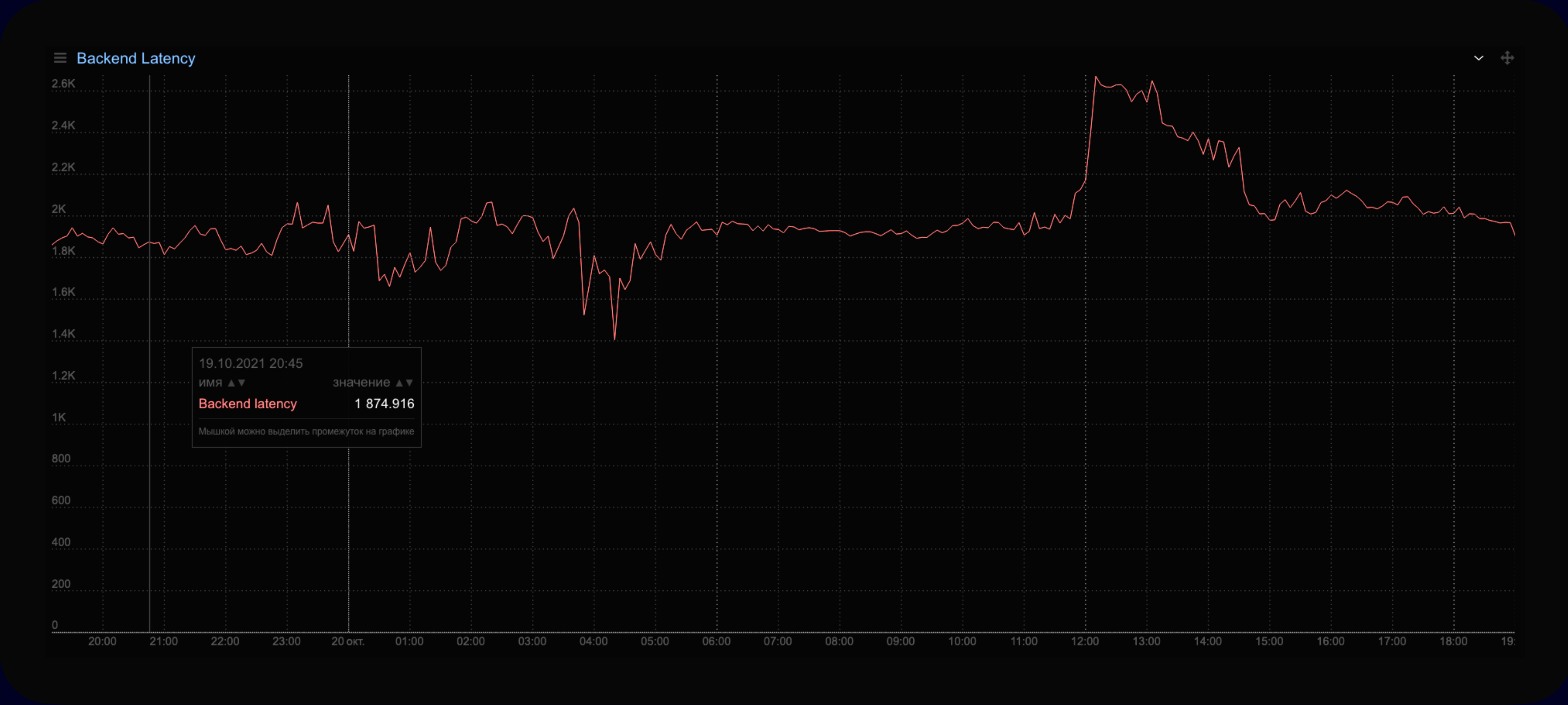
# Метрики



# Раскладываем задержку по компонентам

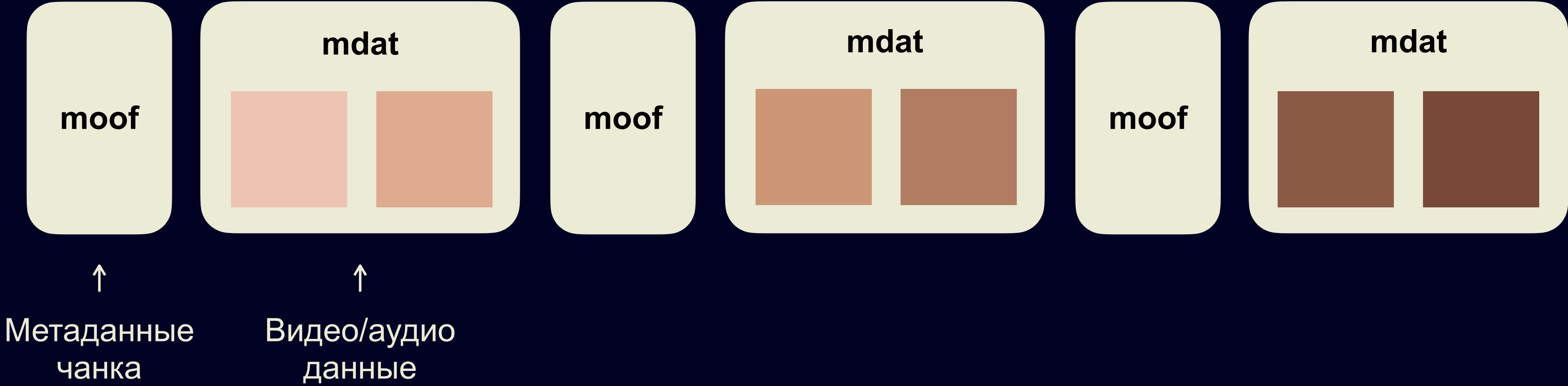


# Итоговая задержка в CDN



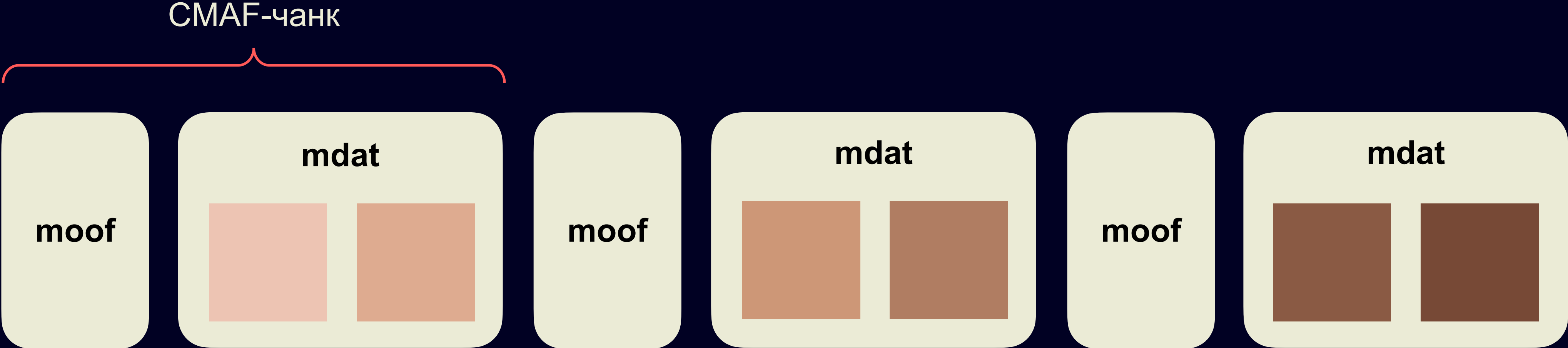
Как показать CMAF  
сегменты в плеере?

# Как устроен CMAF-сегмент?





# Как устроен CMAF-сегмент?



Нужно уметь читать  
тело ответа  
в процессе загрузки

# XMLHttpRequest

XMLHttpRequest 

# Fetch API

# Fetch API + ReadableStream API

# ReadableStream API

Usage

% of all users



Global

94.61%

Current aligned

Usage relative

Date relative

Filtered

All



IE	Edge *	Firefox	Chrome	Safari	Opera	Safari on iOS *	Opera Mini *	Android Browser *
	12-13	2-64	4-42	3.1-10	10-29	3.2-10.2		
6-10	14-93	65-92	43-94	10.1-14.1	30-79	10.3-14.8		2.1-4.4.4
11	94	93	95	15	80	15	all	94
		94-95	96-98	TP				

```
const response = await fetch('segment.m4s');
const reader = response.body.getReader();

async function pump() {
  const { value: chunk, done } = await reader.read();
  if (done) {
    return;
  }

  pushToSourceBuffer(chunk);
  return pump();
}

await pump();
```



```
const response = await fetch('segment.m4s'); ← Получили
const reader = response.body.getReader();      заголовок ответа

async function pump() {
  const { value: chunk, done } = await reader.read();
  if (done) {
    return;
  }

  pushToSourceBuffer(chunk);
  return pump();
}

await pump();
```

```
const response = await fetch('segment.m4s');
const reader = response.body.getReader();

async function pump() {
  const { value: chunk, done } = await reader.read();
  if (done) {
    return;
  }

  pushToSourceBuffer(chunk);
  return pump();
}
```

```
await pump(); ← Запустили чтение
                через ReadableStream
```

```
const response = await fetch('segment.m4s');
```

```
const reader = response.body.getReader();
```

```
async function pump() {
```

```
  const { value: chunk, done } = await reader.read();
```

```
  if (done) {
```

```
    return;
```

```
  }
```

```
  pushToSourceBuffer(chunk);
```

```
  return pump();
```

```
}
```

```
await pump();
```



Прочитали  
чанк данных

```
const response = await fetch('segment.m4s');
const reader = response.body.getReader();

async function pump() {
  const { value: chunk, done } = await reader.read();
  if (done) {
    return; ← Тело ответа загружено полностью
              (сервер отдал чанк нулевой длины)
  }

  pushToSourceBuffer(chunk);
  return pump();
}

await pump();
```

```
const response = await fetch('segment.m4s');
const reader = response.body.getReader();

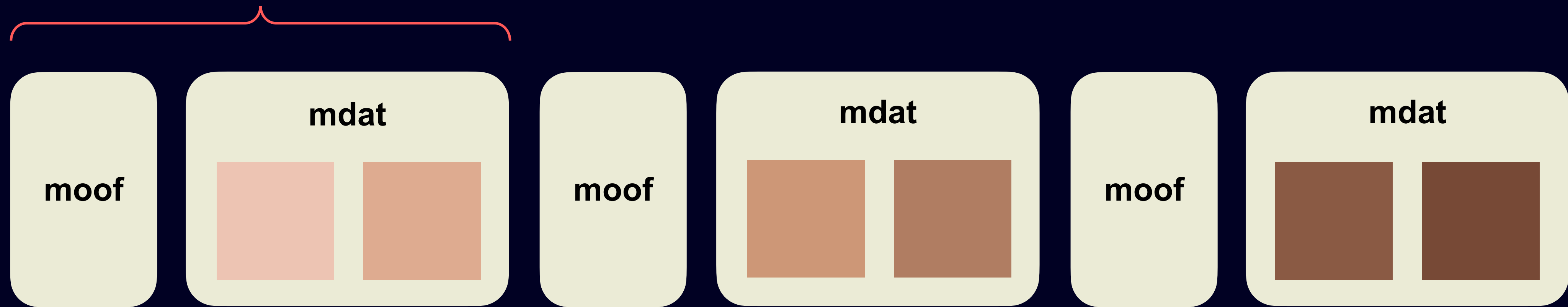
async function pump() {
  const { value: chunk, done } = await reader.read();
  if (done) {
    return;
  }

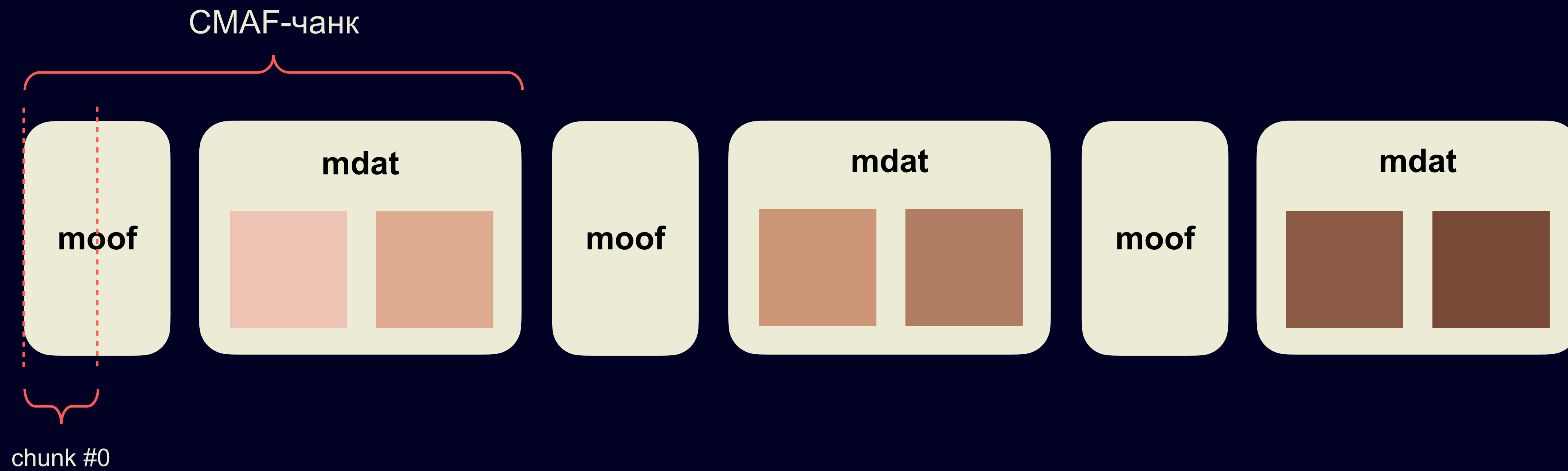
  pushToSourceBuffer(chunk); ← Обрабатываем
  return pump();             чанк сегмента
}

await pump();
```

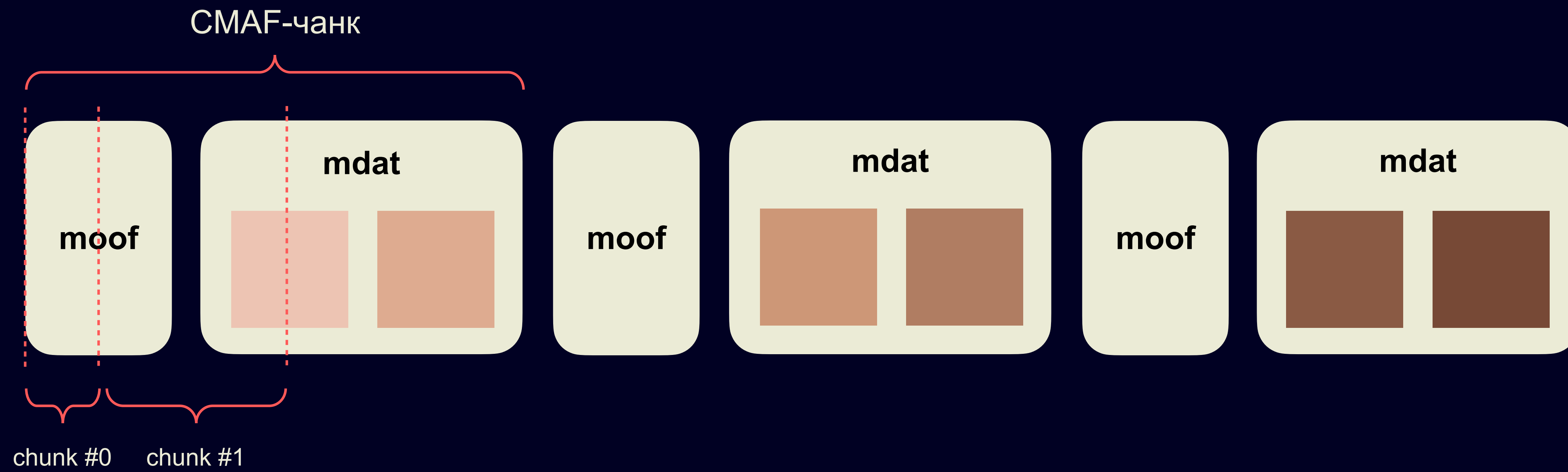
 MEDIA\_ERR\_DECODE

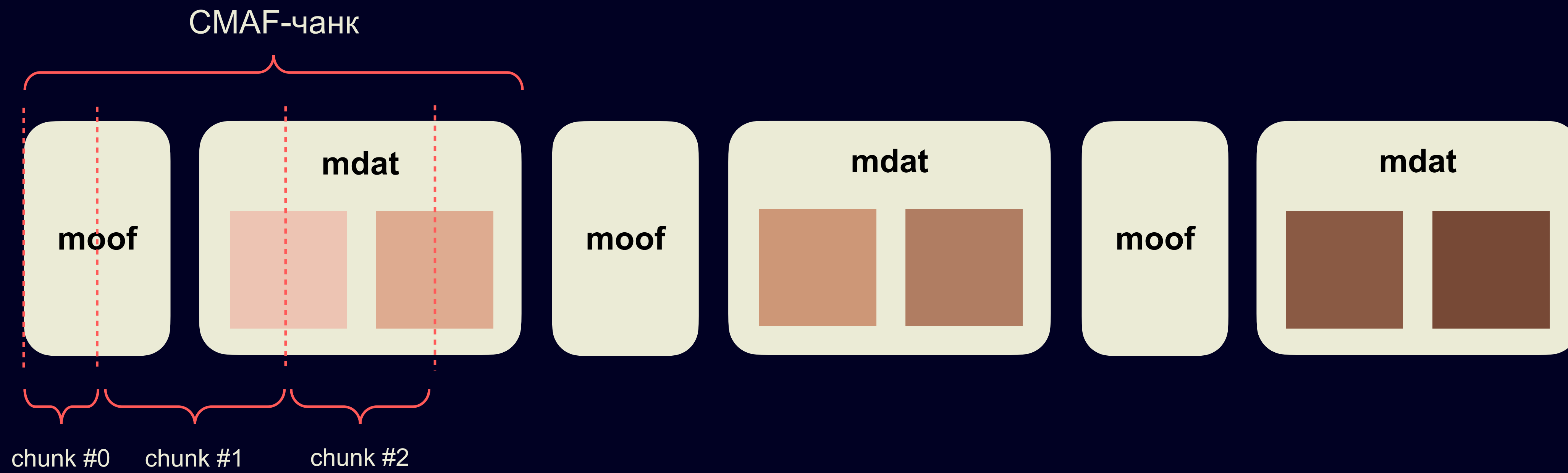
СМАФ-чанк





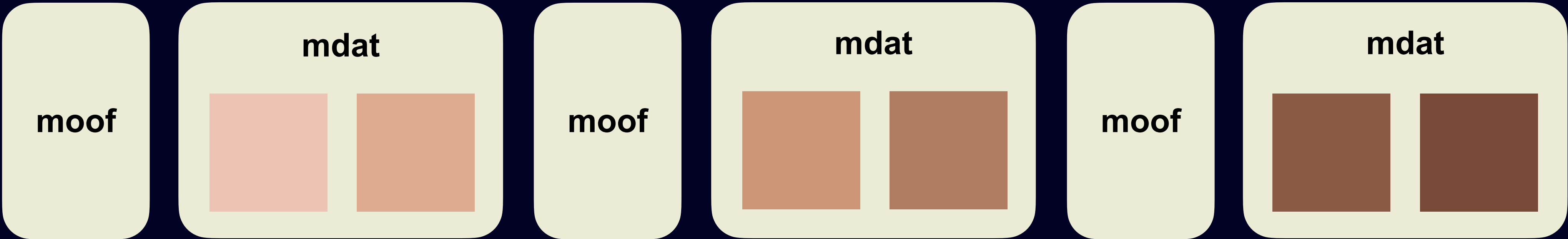




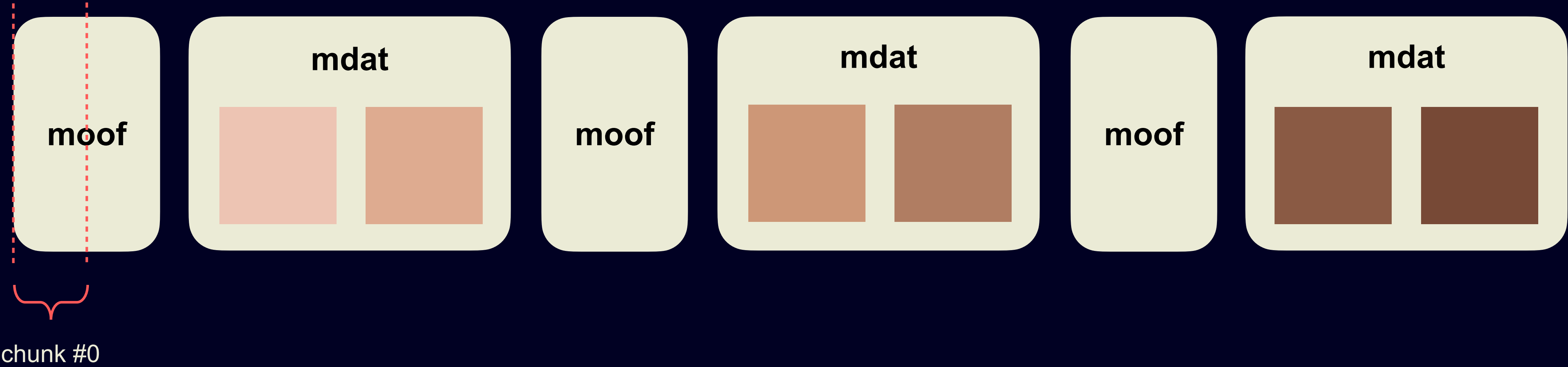


Декодер может обрабатывать  
ТОЛЬКО полные  
mp4-атомы!

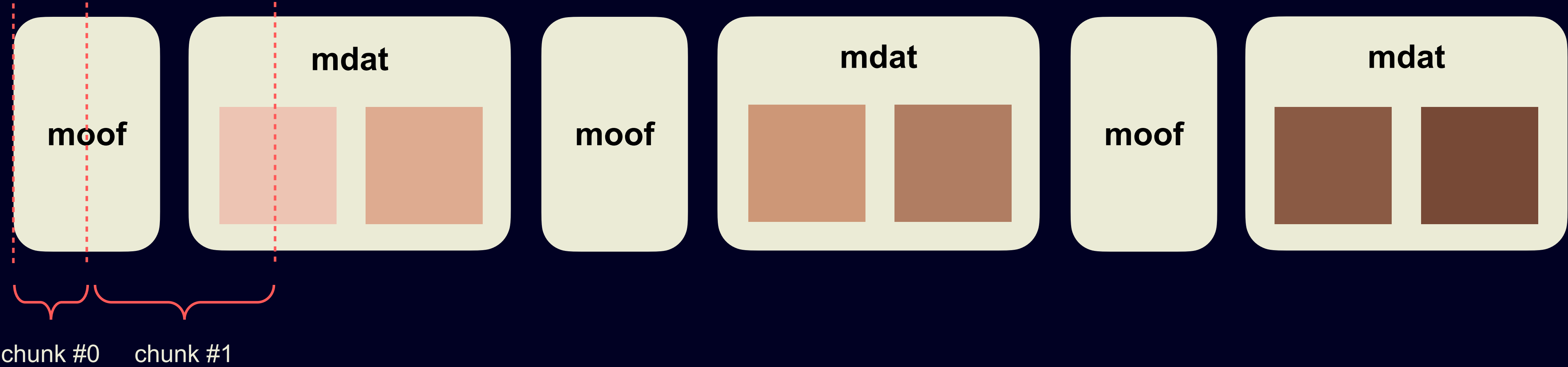
# Парсим tr4 на клиенте



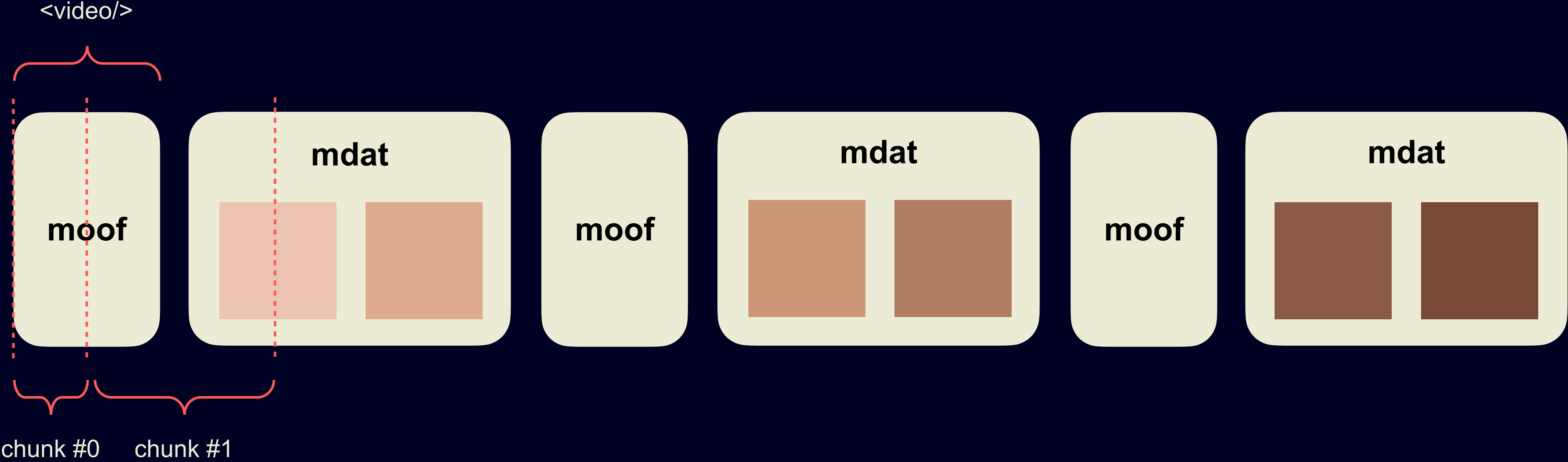
# Парсим tr4 на клиенте



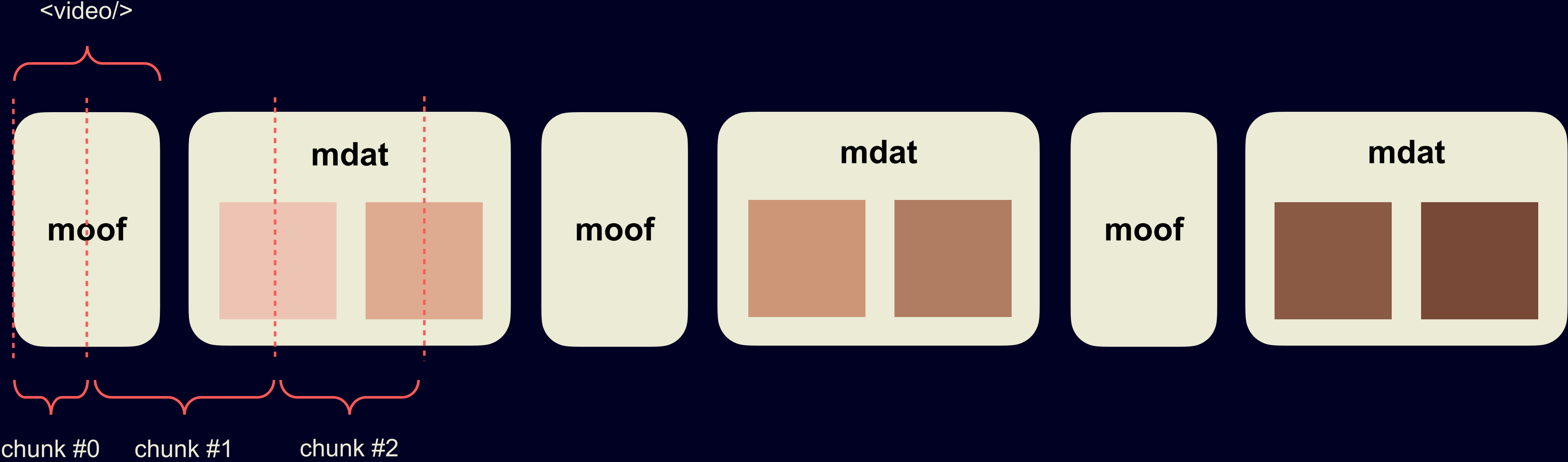
# Парсим tr4 на клиенте



# Парсим mp4 на клиенте

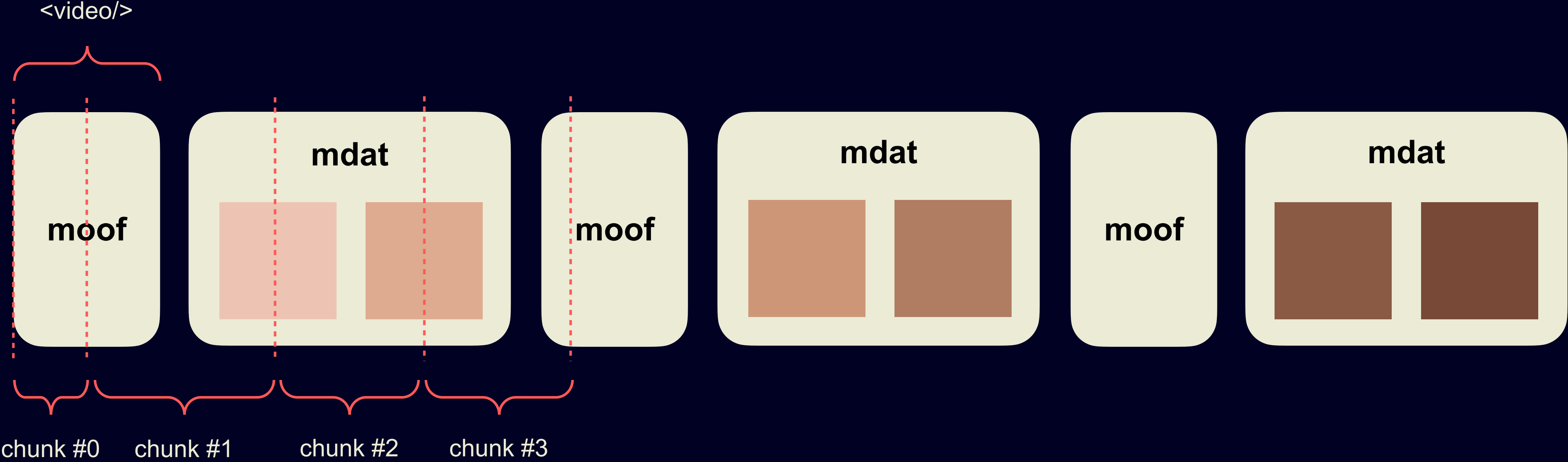


# Парсим mp4 на клиенте

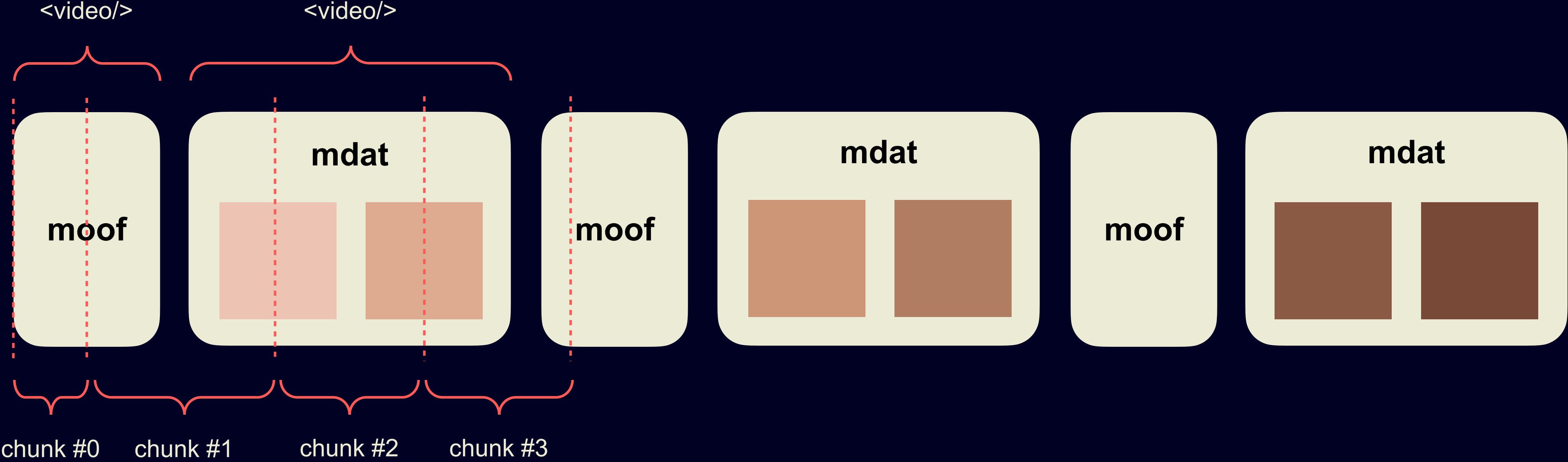




# Парсим mp4 на клиенте



# Парсим mp4 на клиенте



# MP4Box.js



[github.com/gpac/mp4box.js/](https://github.com/gpac/mp4box.js/)

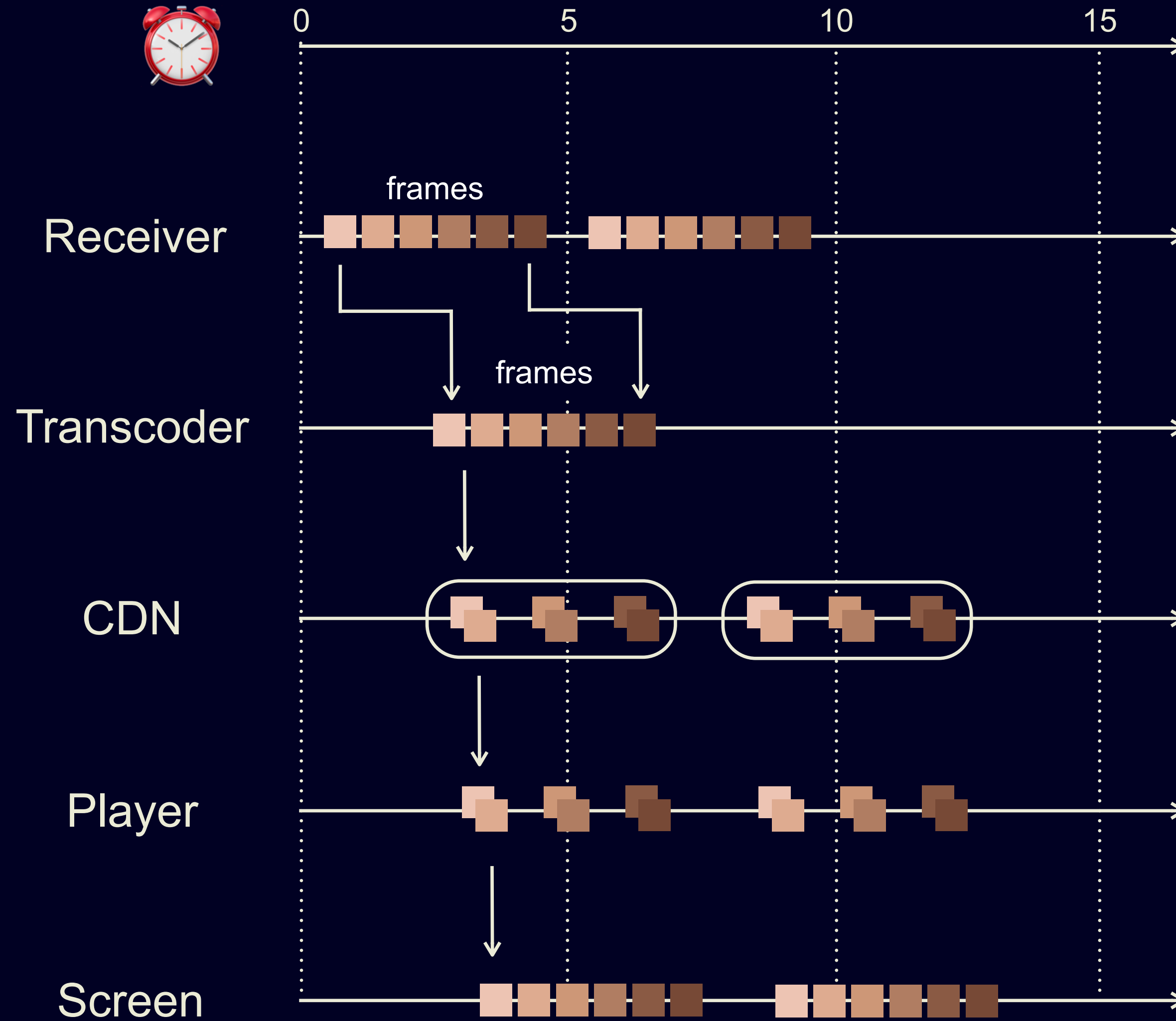
или самописный парсер

Теперь можно делать  
low latency!

Target latency —  
целевое время отставания

Стартовая позиция на таймлайне  
вычисляется как отступ  
на *target latency* от правой границы  
таймлайна

Target latency  $\approx$   
бюджет задержки (10 сек)



Буфер больше не  
накапливаем

Растет вероятность  
буферизации!

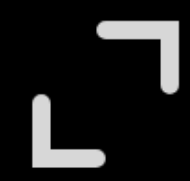


Бюджет задержки сошелся!  
Катим?

Качество Авто 144p ▶

Скорость Обычная ▶

Сообщить о проблеме



10 из 10

КТО ВИНОВАТ,  
И ЧТО ДЕЛАТЬ?

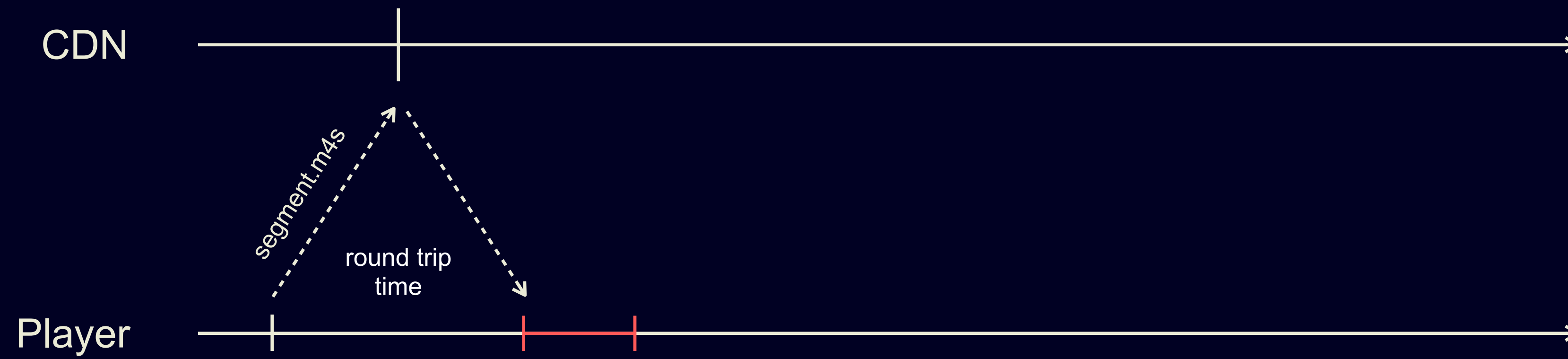
# ABR – Adaptive BitRate

Компонент и набор алгоритмов, отвечающий за выбор видео/аудио дорожки в адаптивном стриминге

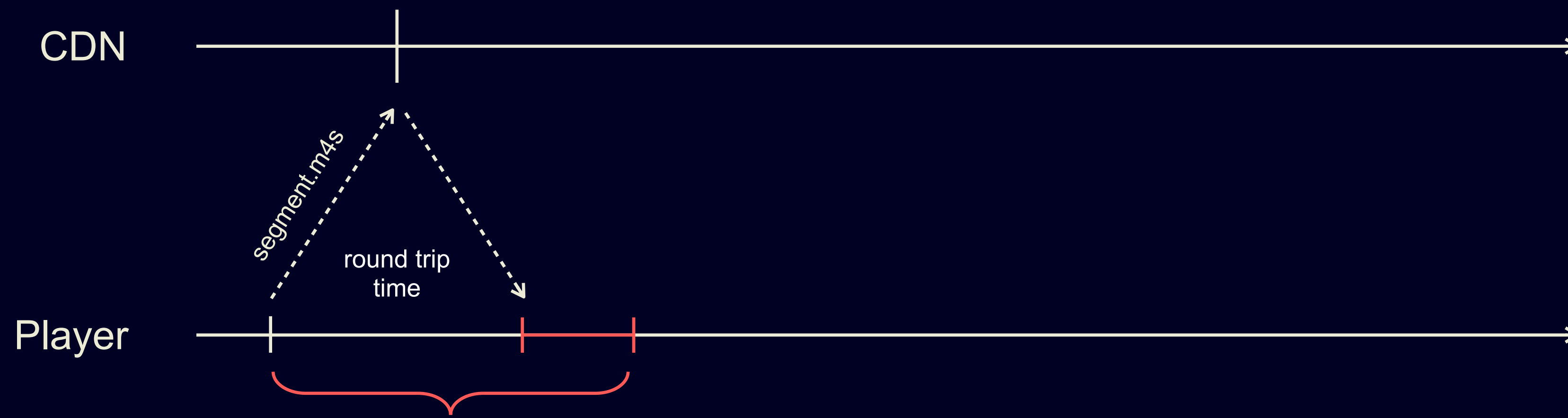
Как устроен  
AVR-менеджер?









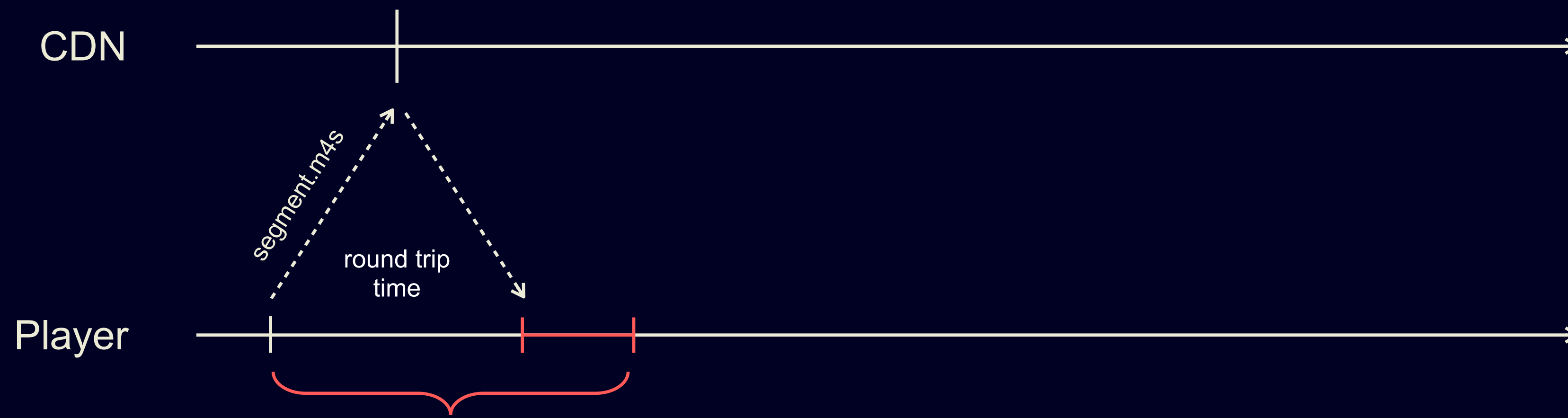


Загружено байт

$b_0$

Время загрузки

$t_0$



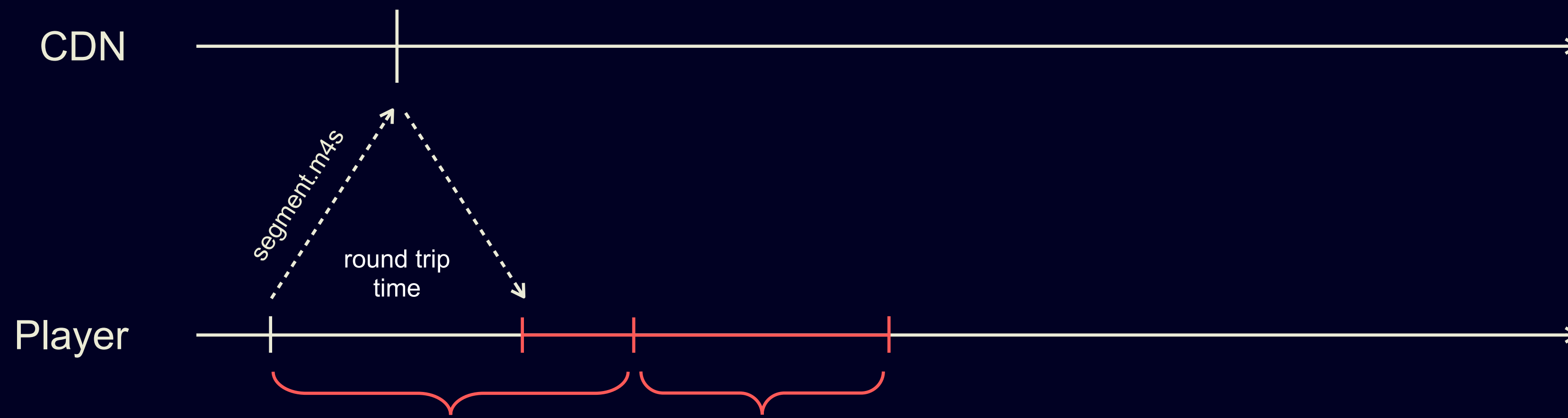
Загружено байт

$b_0$

Время загрузки

$t_0$

$rtt \ll t_0$



Загружено байт

$b_0$

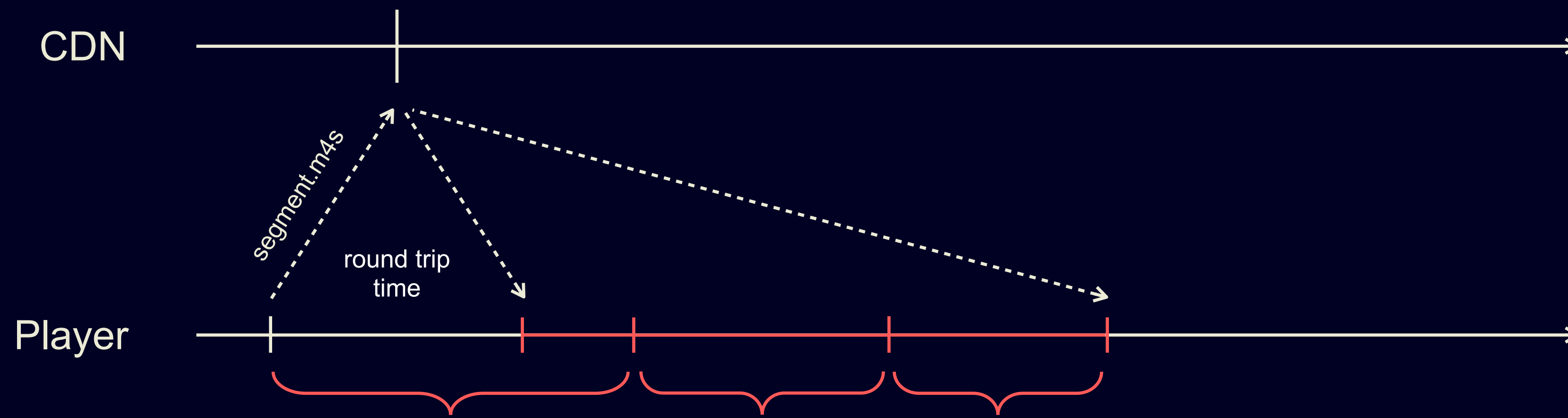
$b_1$

Время загрузки

$t_0$

$t_1$

$$rtt \ll t_0$$



Загружено байт

$b_0$

$b_1$

$b_2$

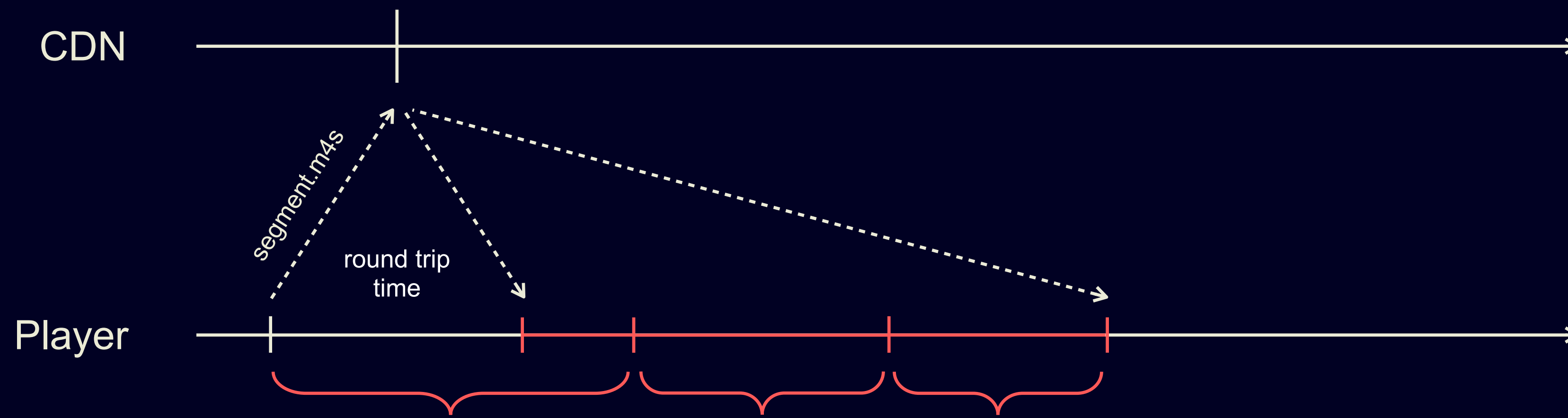
Время загрузки

$t_0$

$t_1$

$t_2$

$$rtt \ll t_0$$



Загружено байт

Время загрузки

$$\frac{b_0}{t_0} = bw_0 \quad \frac{b_1}{t_1} = bw_1 \quad \frac{b_2}{t_2} = bw_2$$

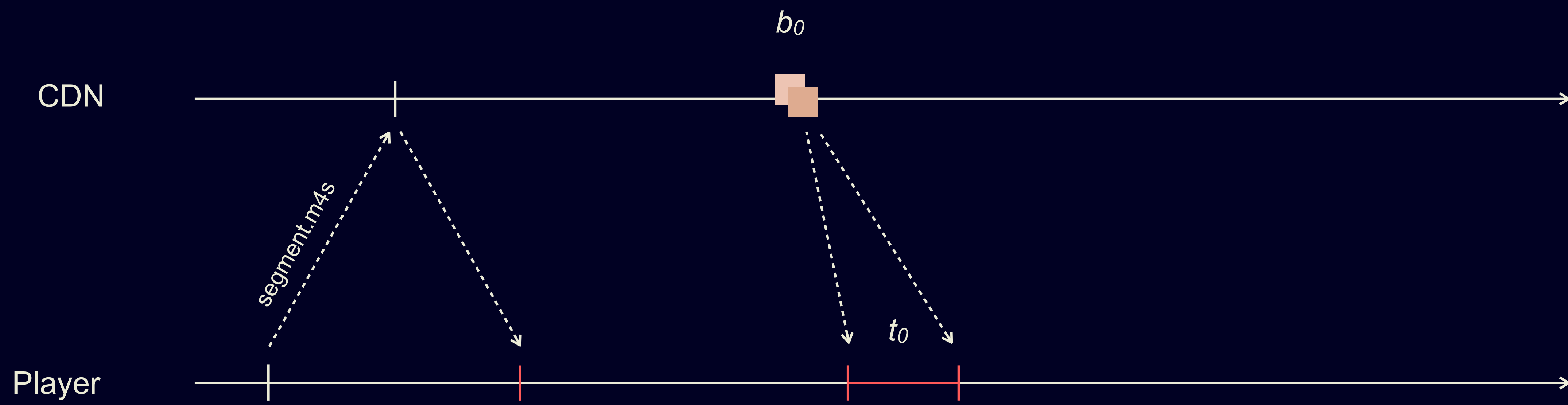
$$rtt \ll t_0$$

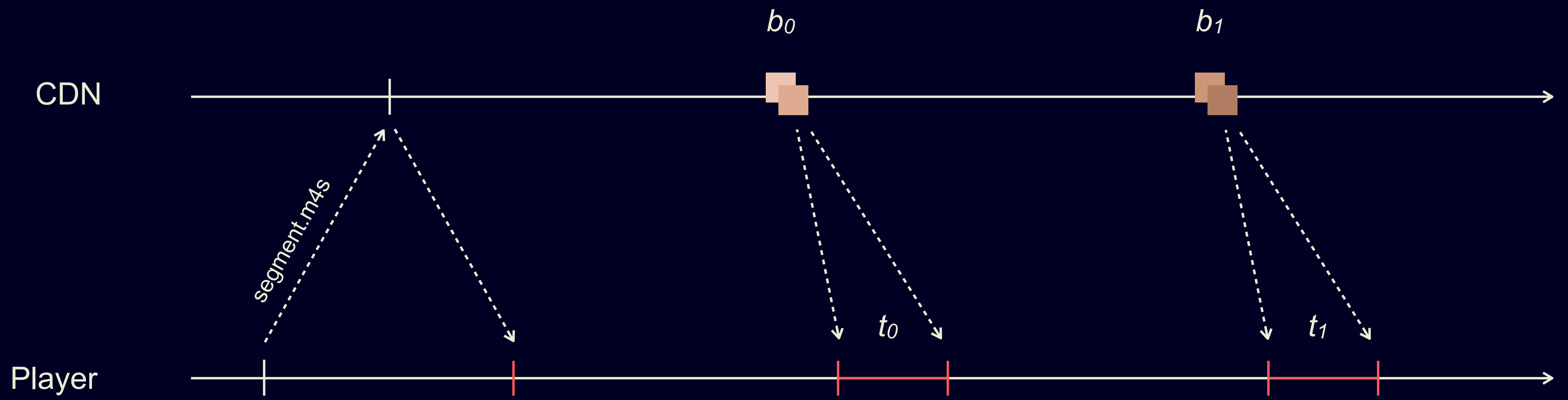
Что происходит  
со СМАФ?

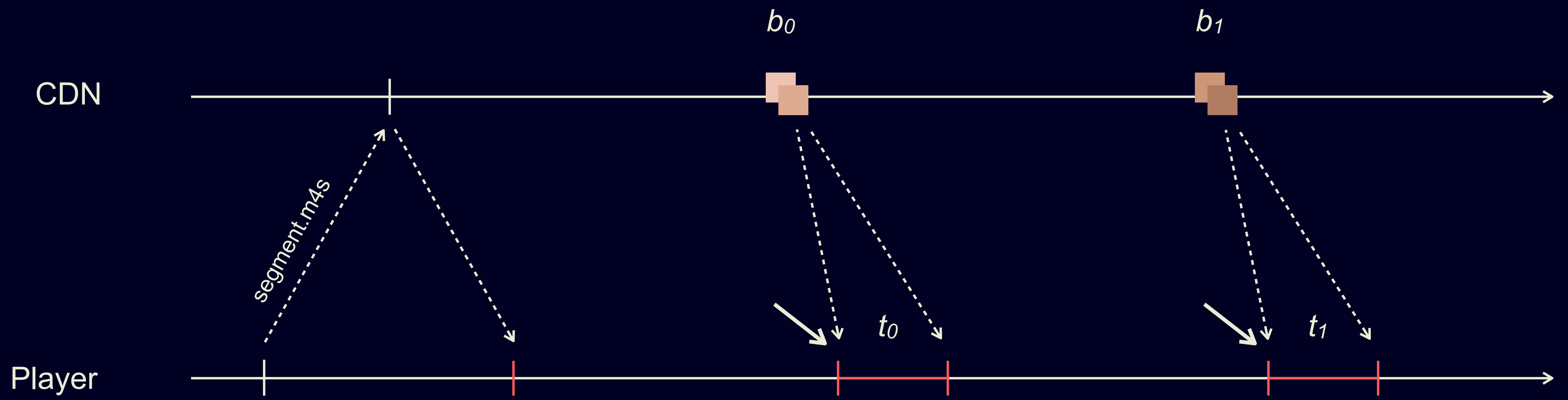


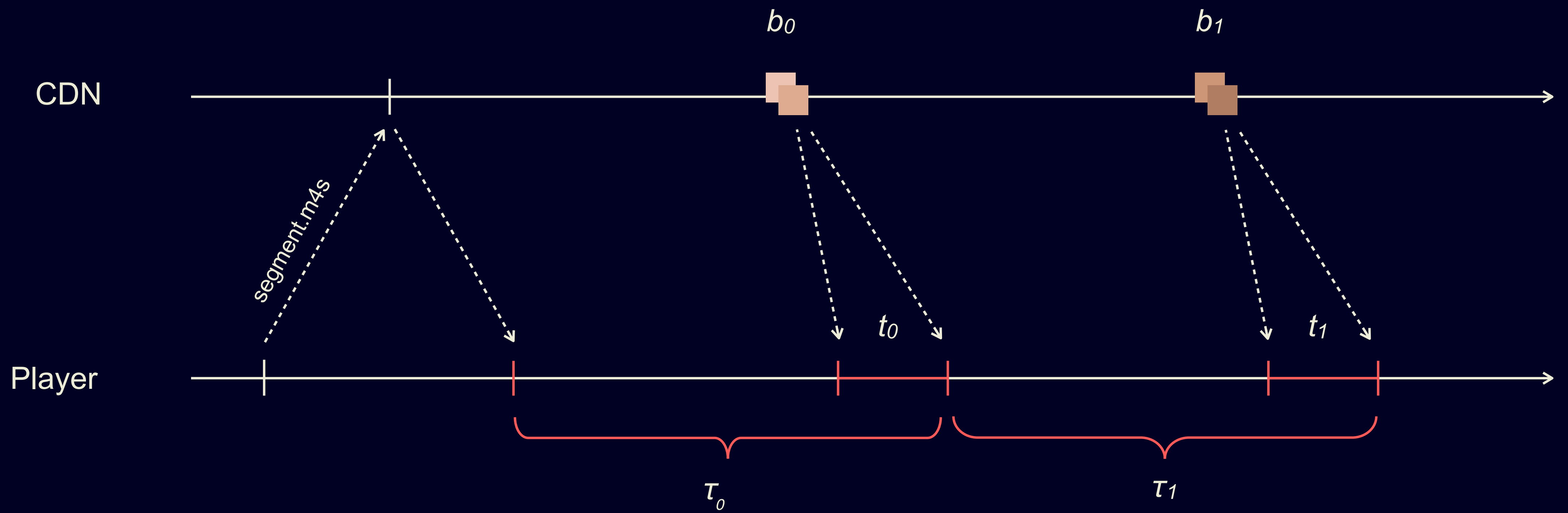


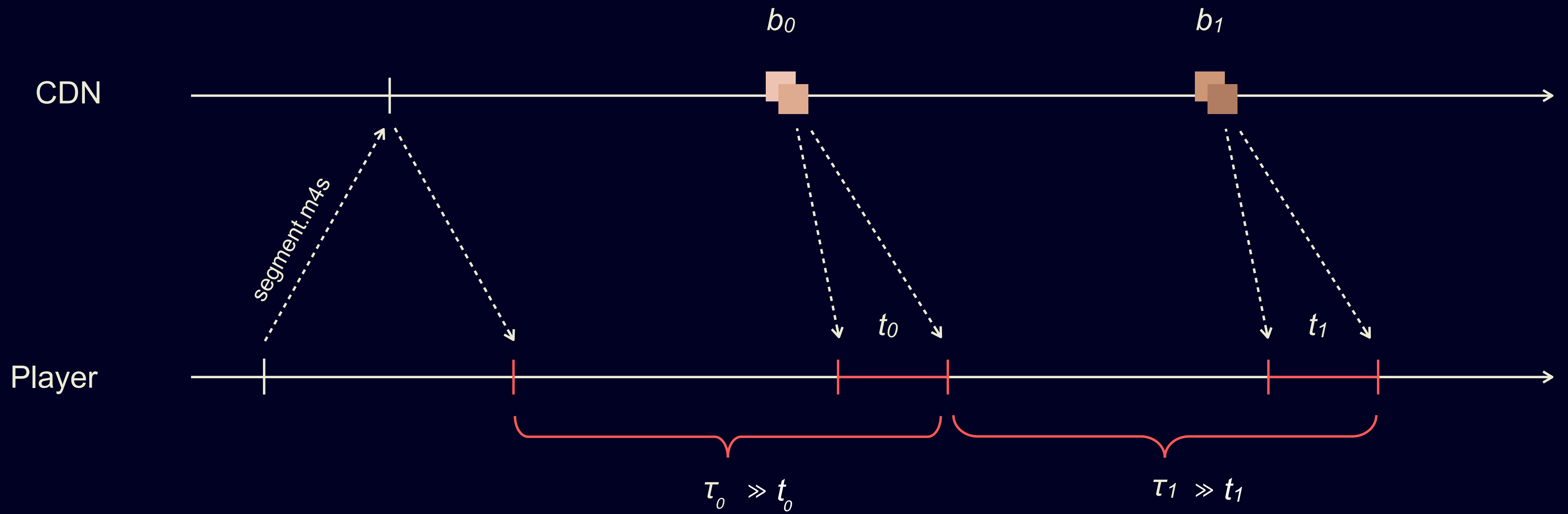












Реальная скорость сети:  $\frac{b}{t}$

Вычисляемая скорость сети:  $\frac{b}{\tau}$

Реальная скорость сети:  $\frac{b}{t}$        $\frac{b}{\tau} \ll \frac{b}{t}$

Вычисляемая скорость сети:  $\frac{b}{\tau}$

Реальная скорость сети:  $\frac{b}{t}$

$$\frac{b}{\tau} \ll \frac{b}{t}$$



Вычисляемая скорость сети:  $\frac{b}{\tau}$

144p



AVR-менеджер нужно  
адаптировать для low latency

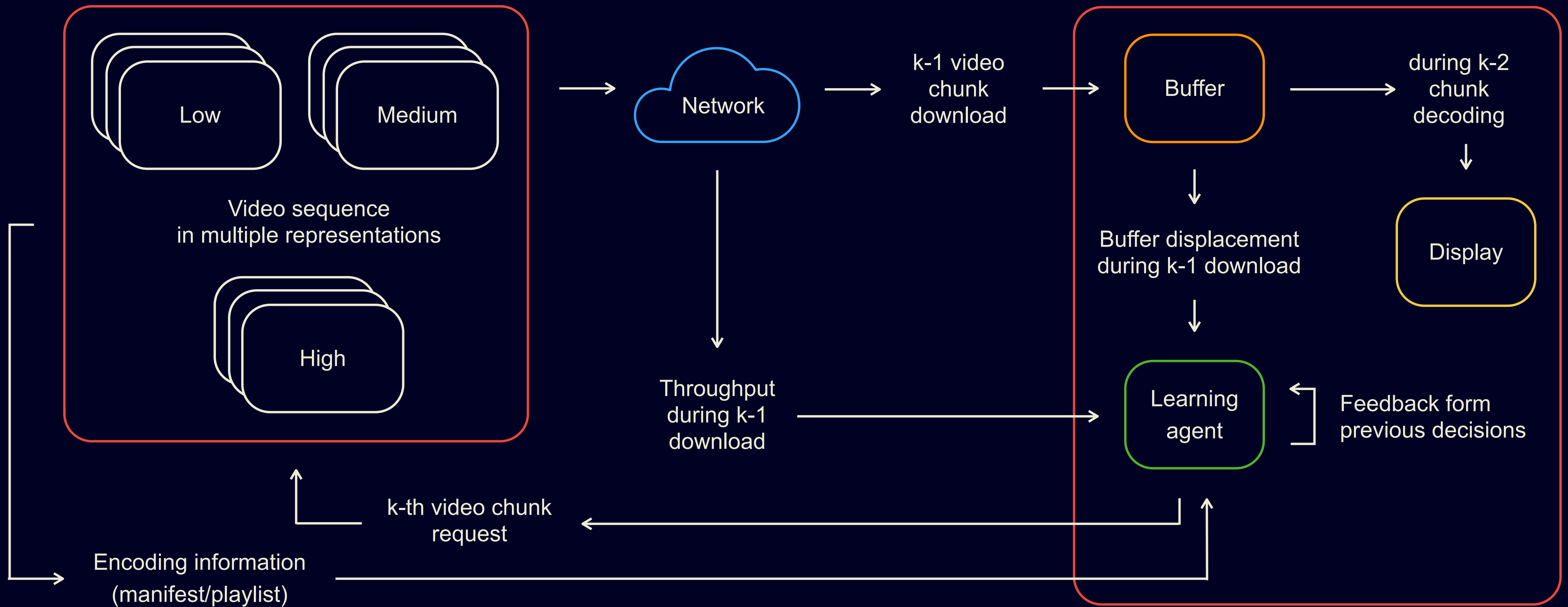
А как решают проблему  
другие сервисы?

ACM MMSys 2020 Grand Challenge on  
**Adaptation Algorithms for Near-Second Latency**

Organized and sponsored by



# Learn2Adapt-LowLatency (L2A-LL)





Имплементация L2A-LL

[github.com/Dash-Industry-Forum/dash.js/pull/3424](https://github.com/Dash-Industry-Forum/dash.js/pull/3424)









Stream <https://cmafref.akamaized.net/cmaf/live-ull/2006350/akambr/out.mpd> Show Options Stop Load



Additional samples can be found in the [Sample Section](#).

- Video Audio
- Buffer Length : 2.993
  - Bitrate Downloading : 6000 kbps
  - Index Downloading : 7 / 7
  - Index playing : 7 / 7
  - Dropped Frames : 5
  - Latency (min|avg|max) : 0.01 | 0.10 | 0.19
  - Download (min|avg|max) : 0.00 | 1.28 | 1.82
  - Ratio (min|avg|max) : 1.10 | 1.56 | 667.33
  - Live Latency: 3.006

Elements Console Sources Network Performance Memory Application Security Lighthouse Requestly

Filter  Invert  Hide data URLs **All** Fetch/XHR JS CSS Img Media Font Doc WS Wasm Manifest Other  Has blocked cookies  Blocked Requests  3rd-party requests

20000 ms	40000 ms	60000 ms	80000 ms	100000 ms	120000 ms	140000 ms	160000 ms	180000 ms	200000 ms	220000 ms	240000 ms	260000 ms
[Network activity visualization]												

Stream <https://cmafref.akamaized.net/cmaf/live-ull/2006350/akambr/out.mpd> Show Options Stop Load



Additional samples can be found in the [Sample Section](#).

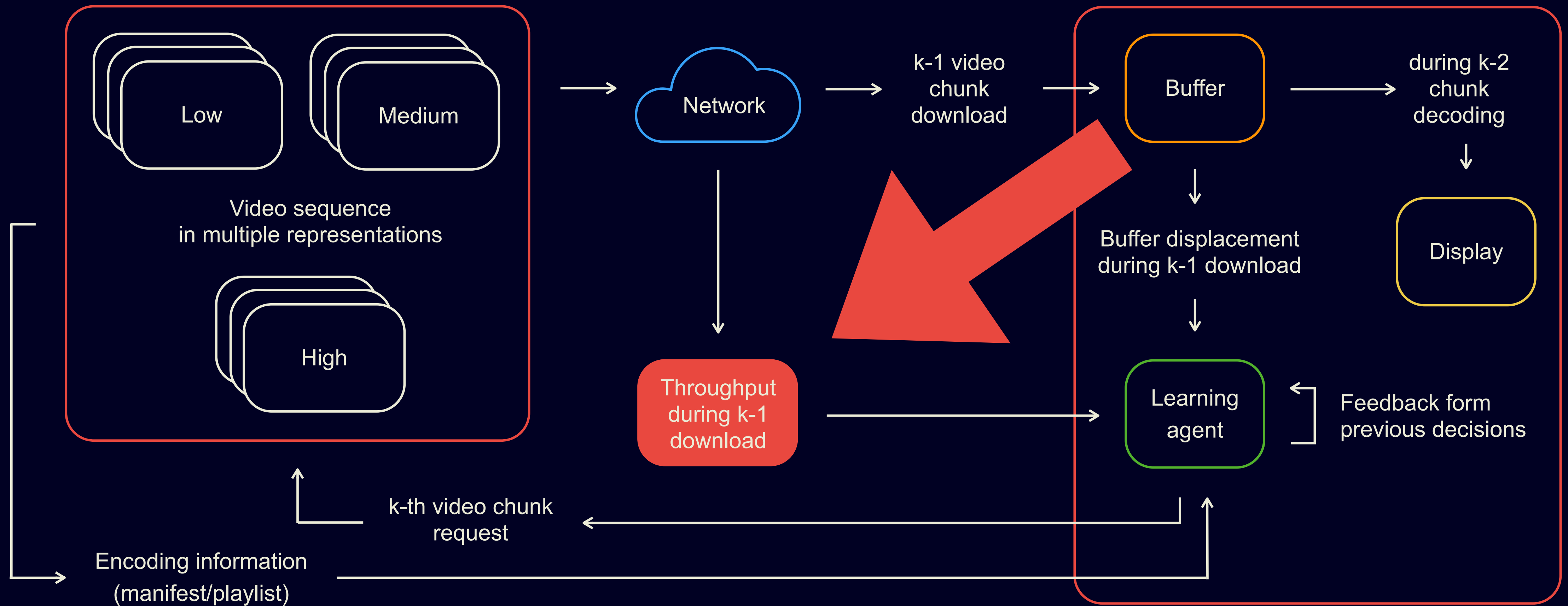
- Video Audio
- Buffer Length : 2.993
  - Bitrate Downloading : 6000 kbps
  - Index Downloading : 7 / 7
  - Index playing : 7 / 7
  - Dropped Frames : 5
  - Latency (min|avg|max) : 0.01 | 0.10 | 0.19
  - Download (min|avg|max) : 0.00 | 1.28 | 1.82
  - Ratio (min|avg|max) : 1.10 | 1.56 | 667.33
  - Live Latency: 3.006

Elements Console Sources Network Performance Memory Application Security Lighthouse Requestly

Filter  Invert  Hide data URLs All Fetch/XHR JS CSS Img Media Font Doc WS Wasm Manifest Other  Has blocked cookies  Blocked Requests  3rd-party requests

20000 ms	40000 ms	60000 ms	80000 ms	100000 ms	120000 ms	140000 ms	160000 ms	180000 ms	200000 ms	220000 ms	240000 ms	260000 ms
[Network activity visualization]												

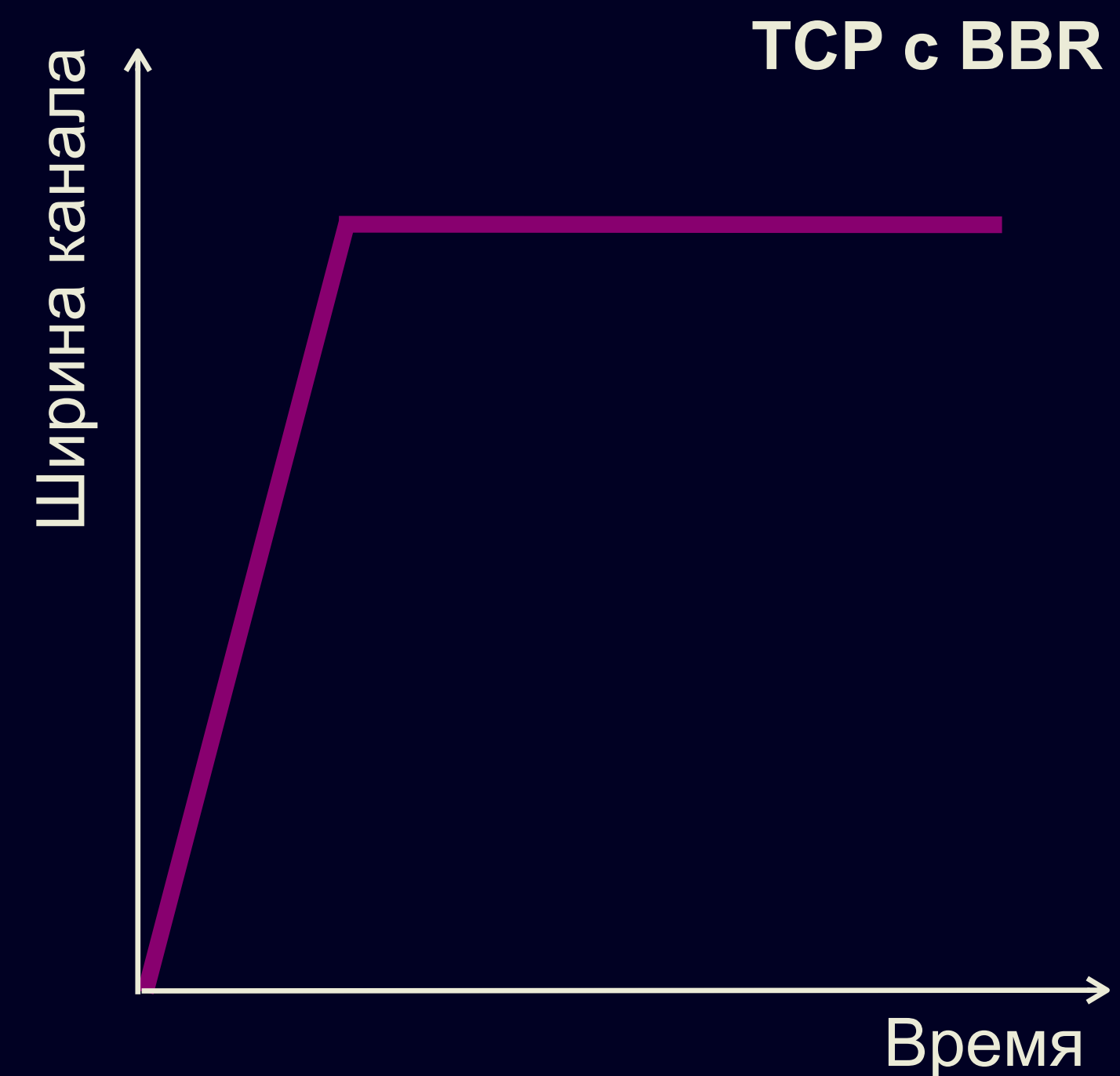
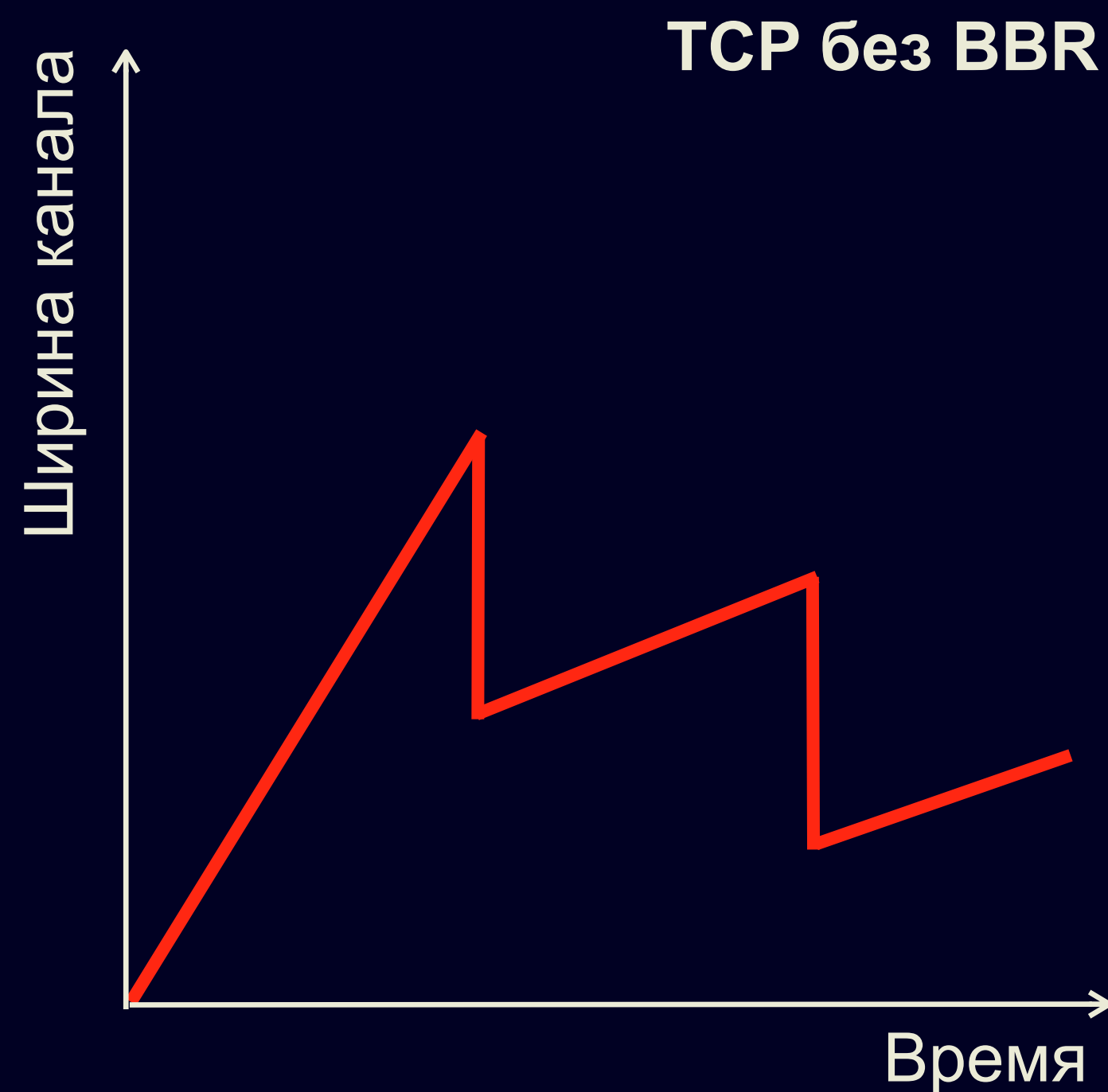
# Learn2Adapt-LowLatency (L2A-LL)



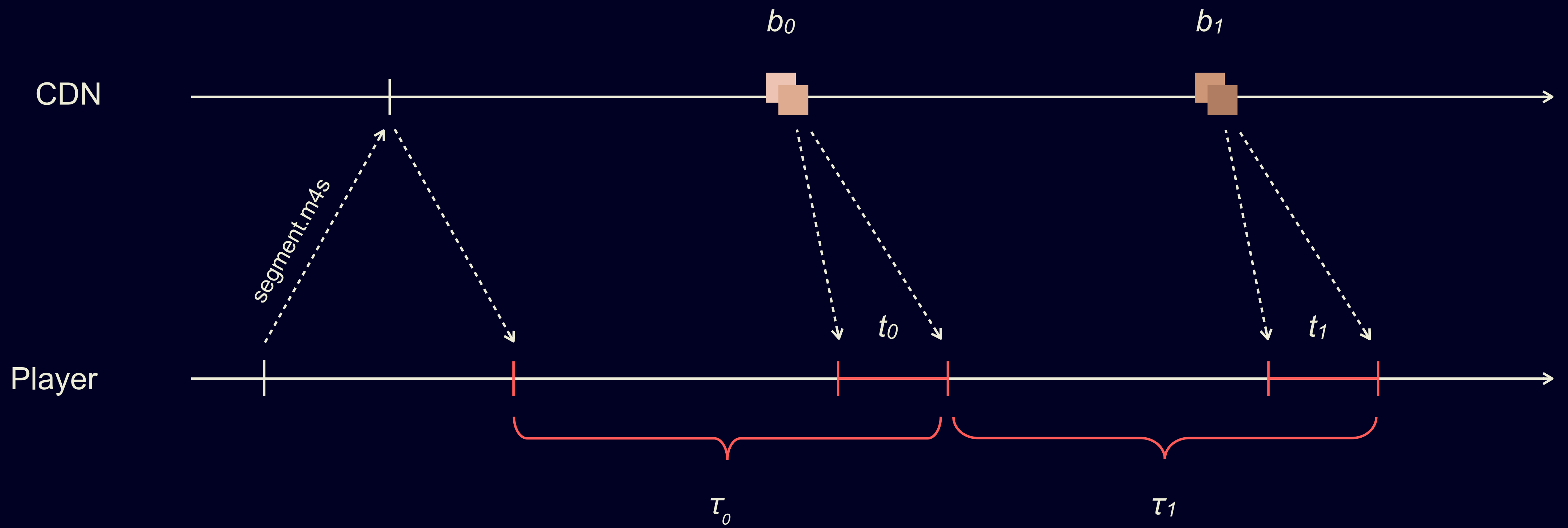
Проблема оценки скорости сети  
никуда не деваается 🤔

Если мы не можем делать что-то  
на клиенте, давайте делать  
это на сервере?

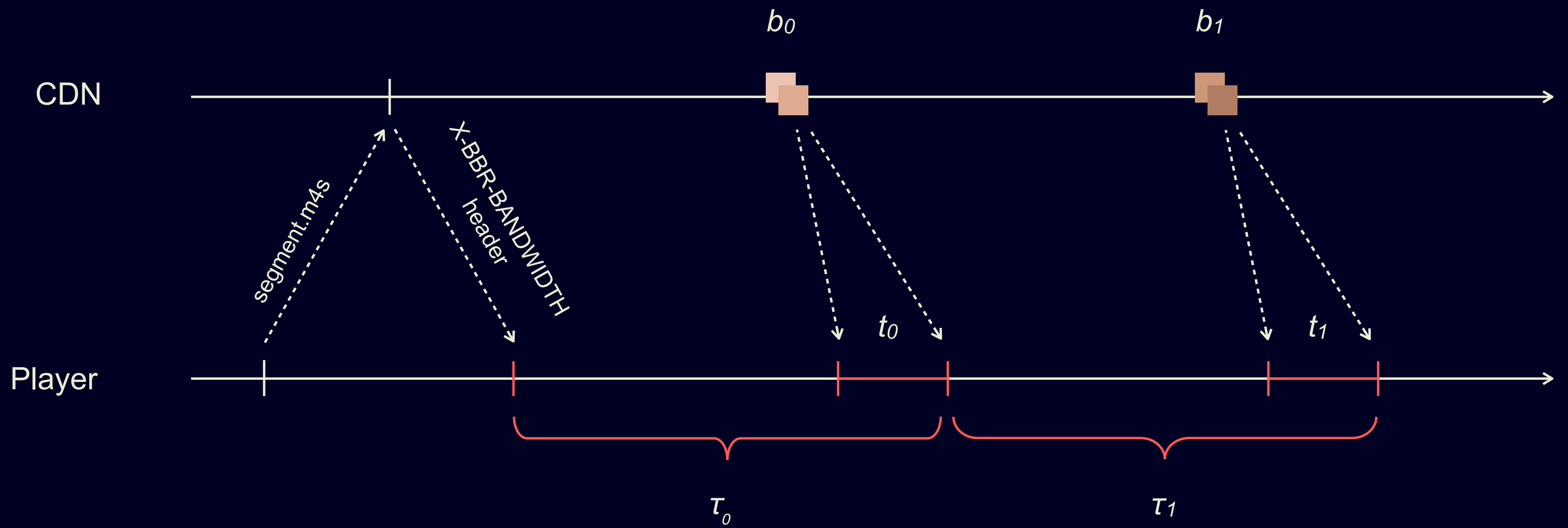
# ТСР BBR: быстрый и простой способ ускорения загрузки страниц

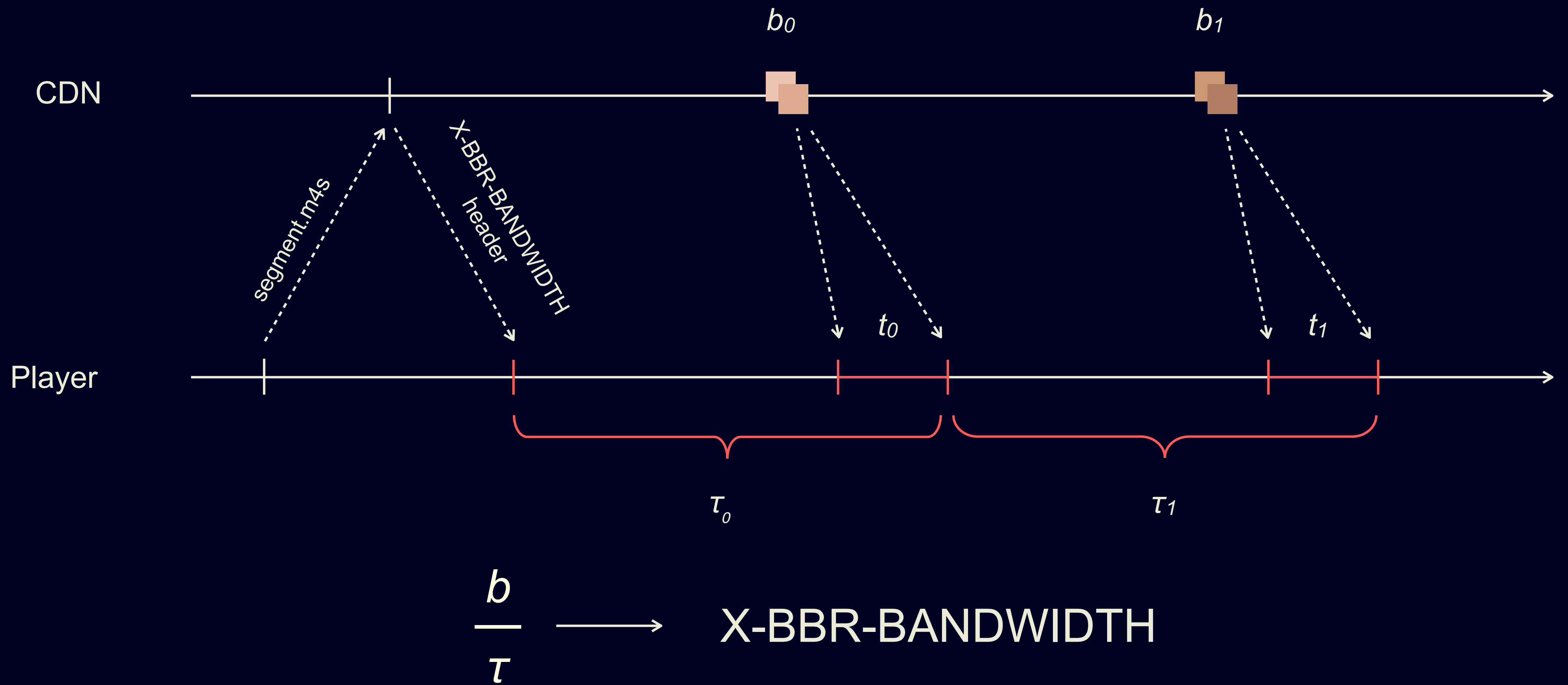


Попробуем?





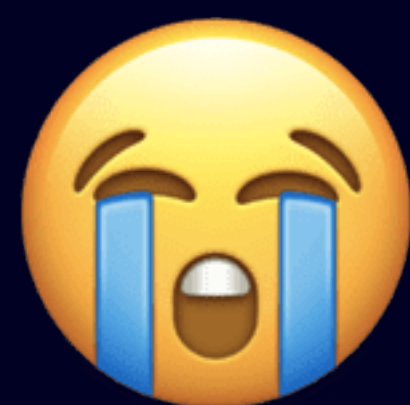




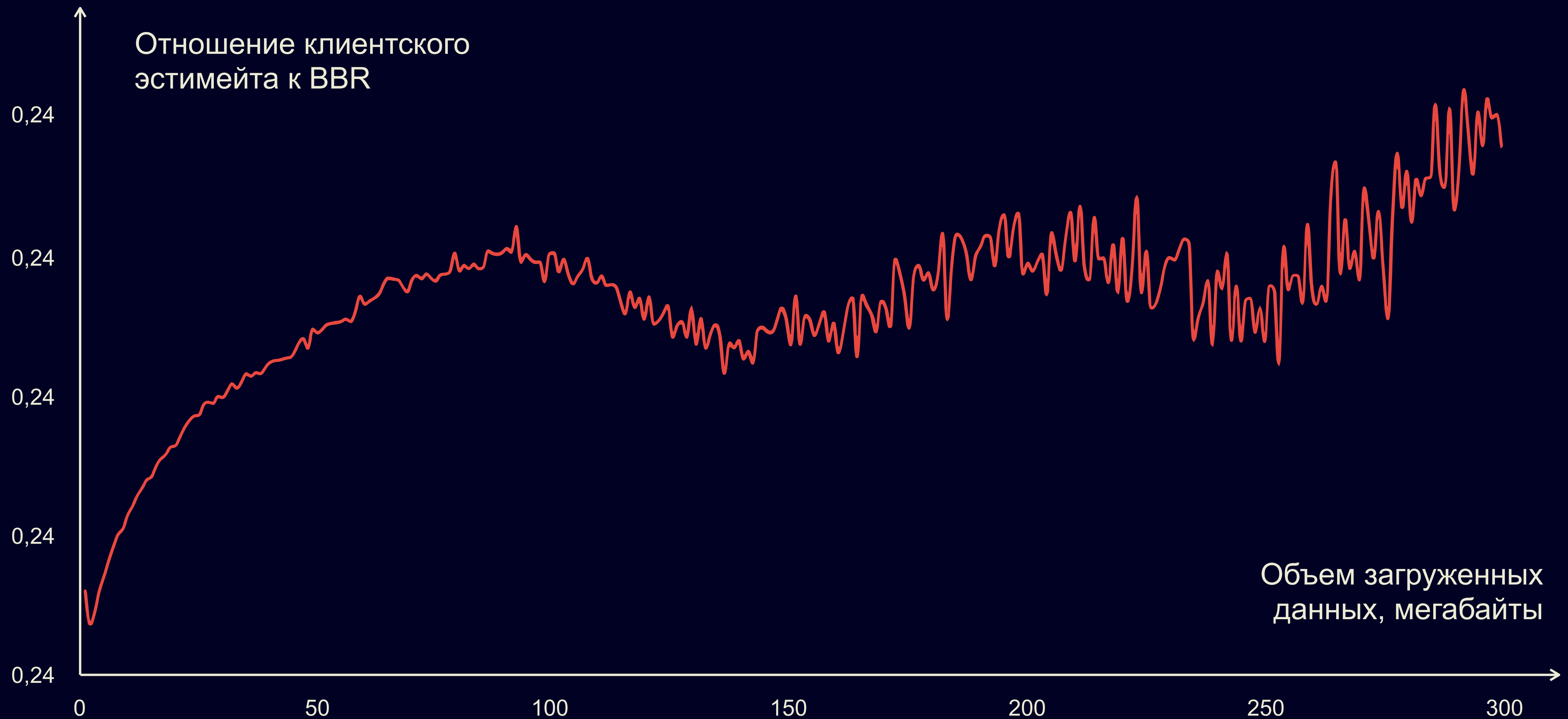
Если на VOD ок —  
то и на LIVE ок

A/B-эксперимент

Сравниваем метрики  
буферизаций



+30% времени  
буферизации



# Выводы

# Выводы

- 1 BBR завышает оценку скорости



# Выводы

- 1 BBR завышает оценку скорости
- 2 BBR медленно сходится

# Выводы

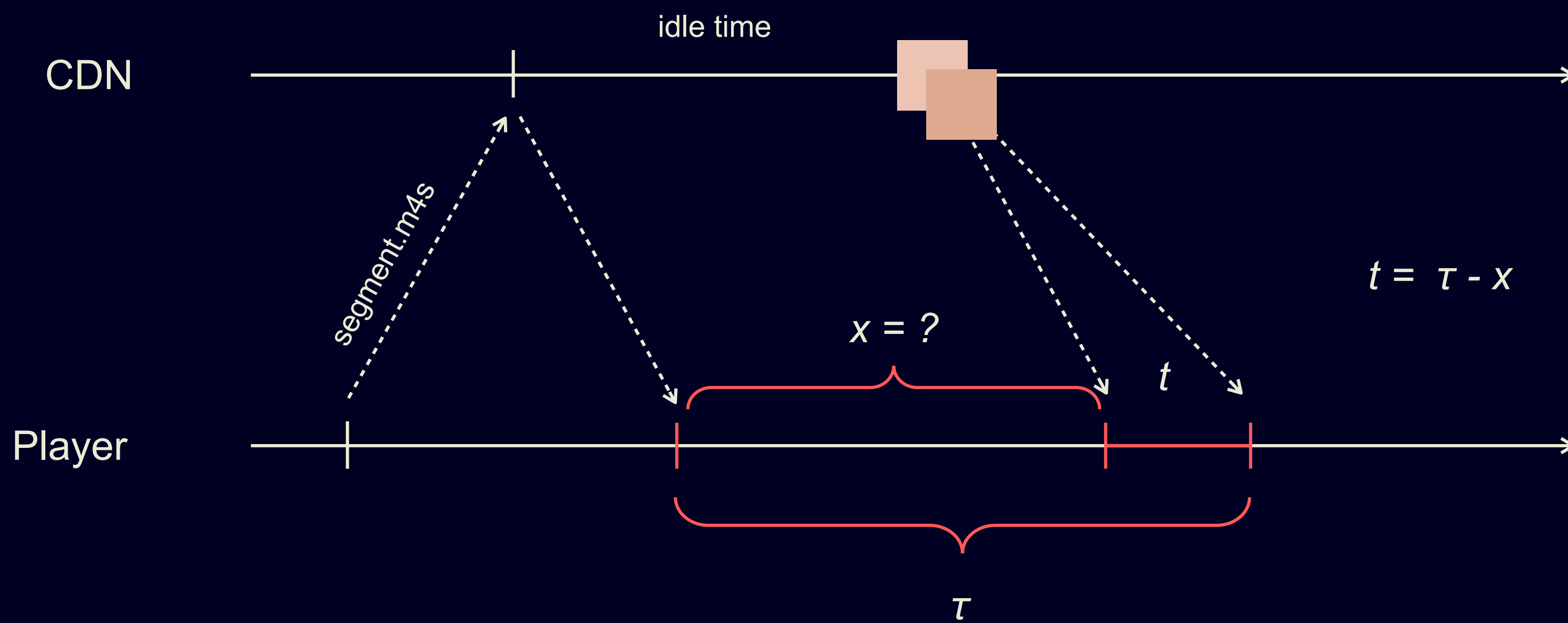
- 1 BBR завышает оценку скорости
- 2 BBR медленно сходится

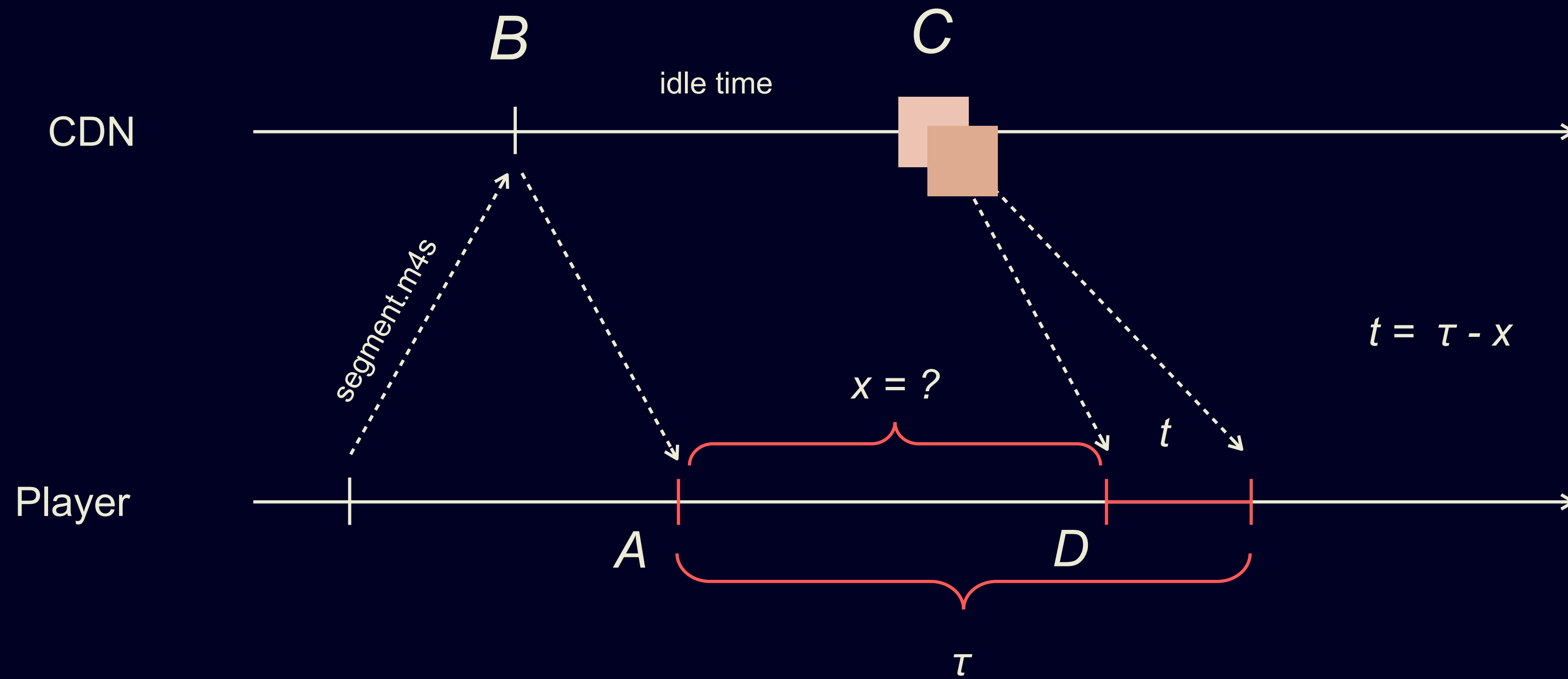
**Вывод:** это решения нам не подходит

И что делать?

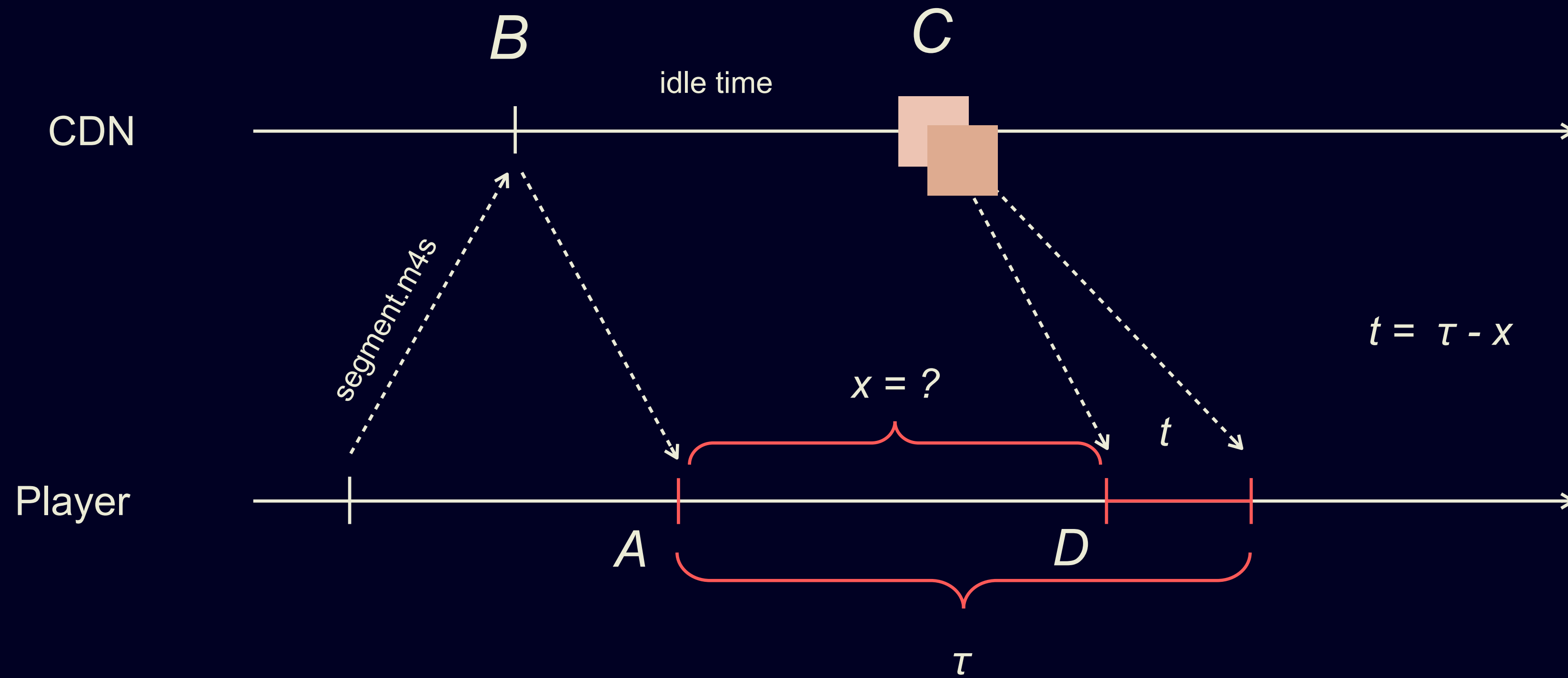






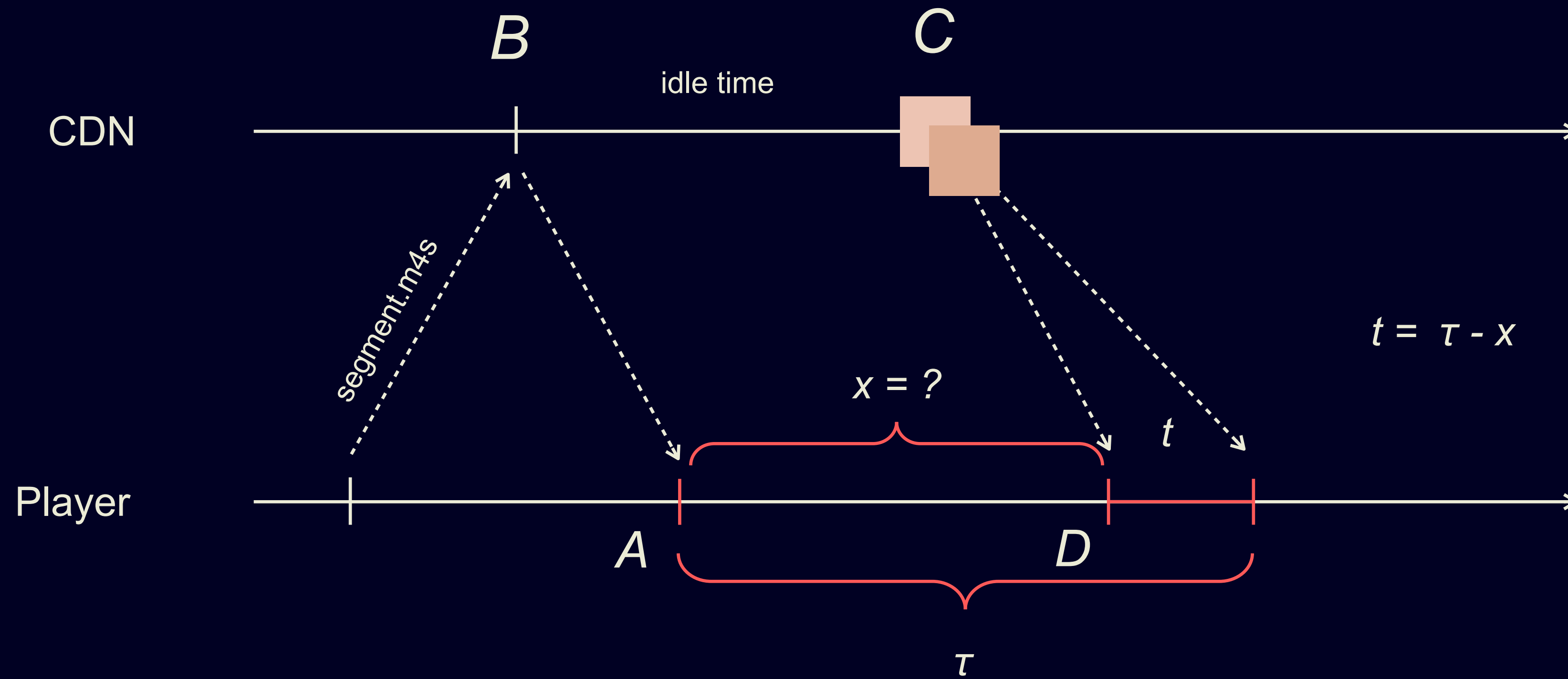


AB || CD



AB || CD

ABCD — параллелограмм

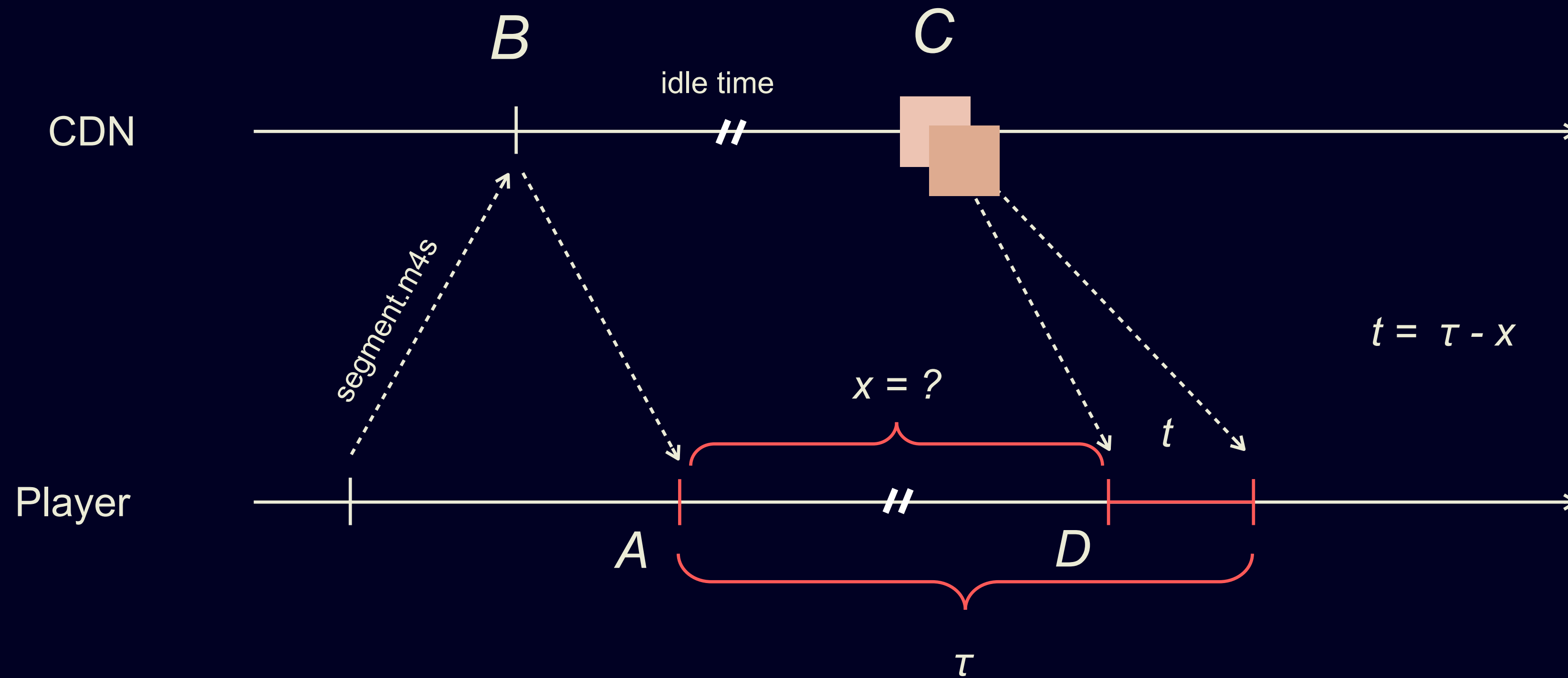




$AB \parallel CD$

ABCD — параллелограмм

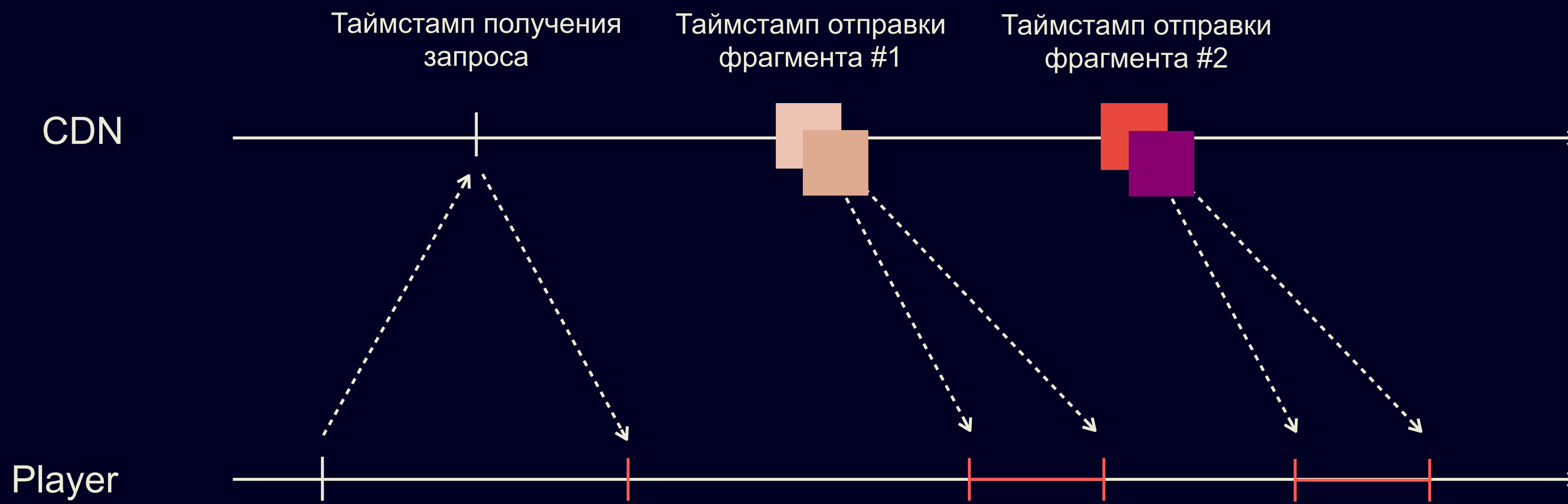
$BC = AD$



$t \approx \tau$  – *idle time*

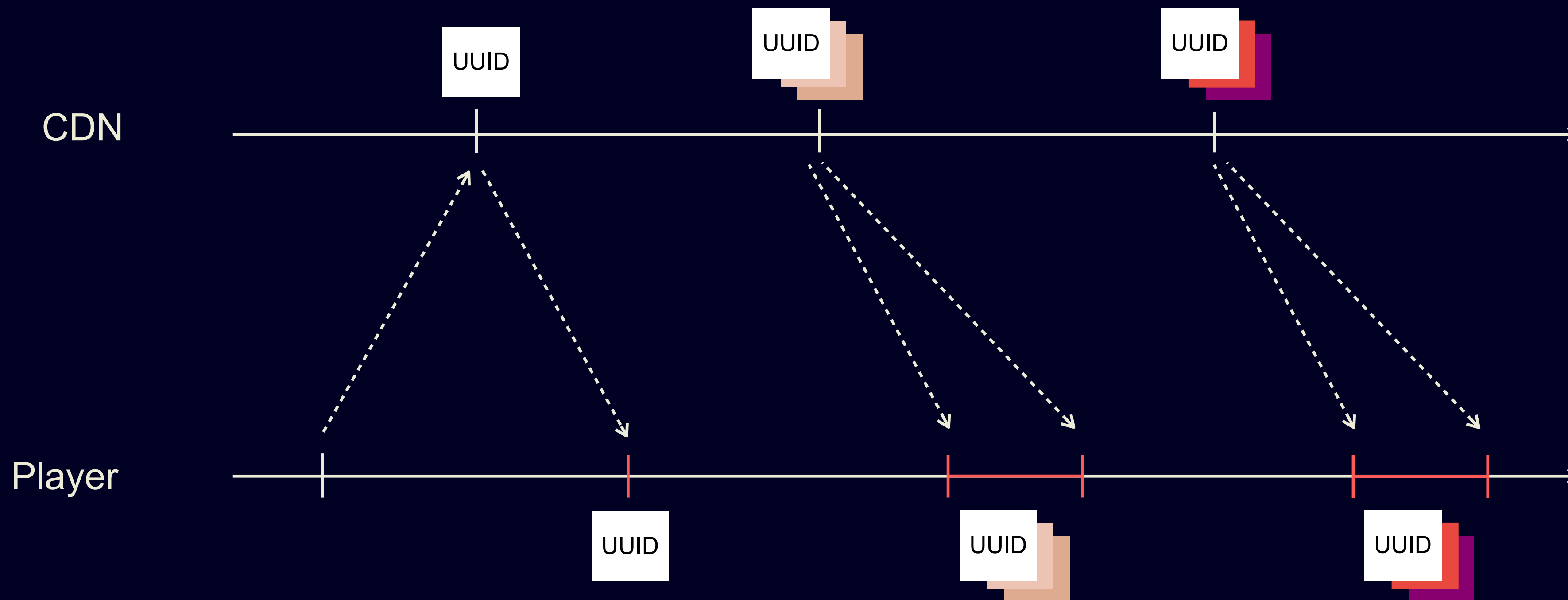
длительность загрузки фрагмента  $\approx$   
длительность получения фрагмента — время «простоя»

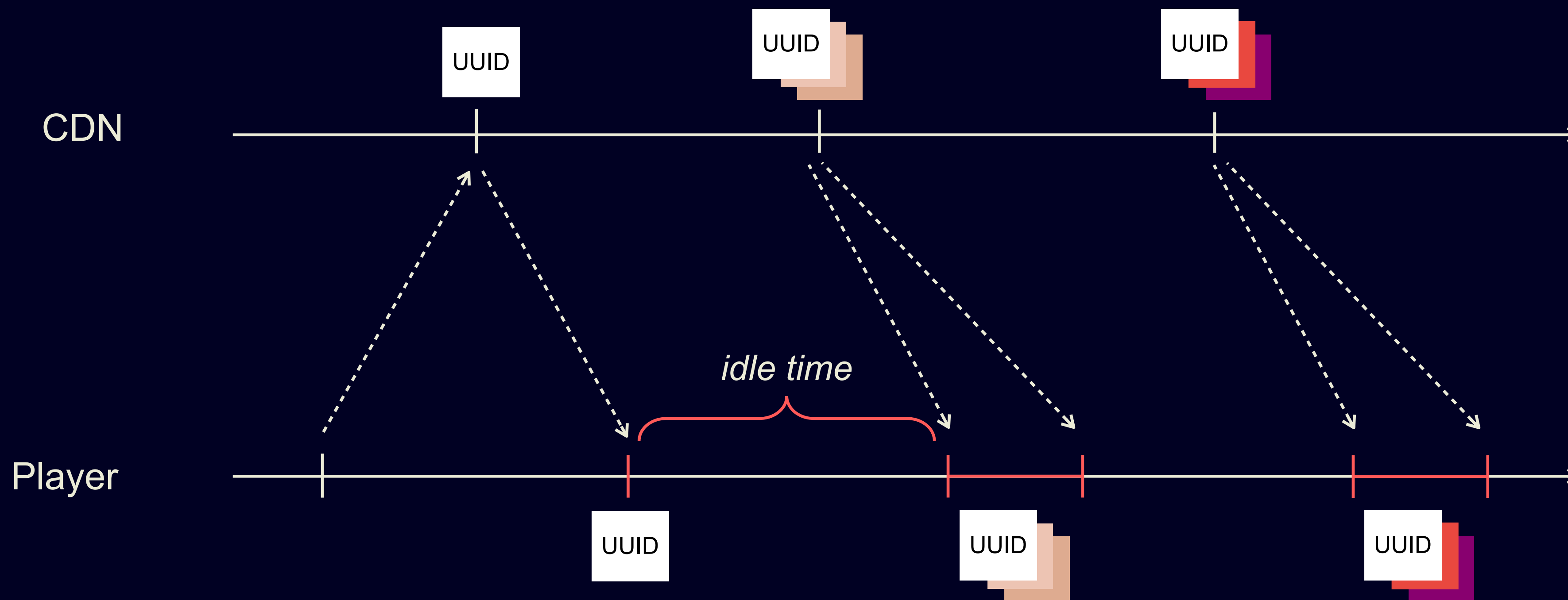
Как получить  
idle time на клиенте?



# UUID-box

- › 16 байт uuid-идентификатор + payload
- › В качестве метаданных в payload запишем серверный таймстамп!





$$\text{idle time} = \text{UUID\#2} - \text{UUID\#1}$$

# Оценка скорости канала

$$bw \approx \frac{b}{\tau - (UUID\#2 - UUID\#1)}$$



# Итоги

- 1 Медиана задержки — 6 секунд  
(при изначальном бюджете — 10 секунд)  
с потенциалом дальнейшего улучшения
- 2 Нет значимых изменений остальных  
продуктовых и технических метрик

# Что мы планируем делать дальше

- 1 Уменьшать задержку за счет дальнейших оптимизаций и "подкрутки" скорости воспроизведения, чтобы нагнать отставание после буферизации
- 2 Улучшить алгоритм оценки скорости сети - сейчас он игнорирует параллельно загружаемые аудио-данные
- 3 Отключать low-latency режим после нескольких прерываний просмотра

# Мы нанимаем!



## Фронтенд

- Разработчик в веб-плеер
- Стажер в веб-плеер

## Android

- Разработчик в android-плеер



**videotech**

**Спасибо**

Дмитрий Кравцов,  
Яндекс.Видеостриминг

Алексей Гусев,  
Яндекс.Дзен

