

From CPU to GPU to ASIC: Our Transcoding Journey

Ilya Mikhaelis Streaming Backend Tech Lead Mayflower







2

Plan

- 1. Introduction: Who we are
- 2. Transcoding pipeline
- 3. Software transcoder
- 4. GPU transcoder
- 5. ASIC transcoder
- 6. Results





🔨 NETINT

Who we are



- 1B+ total visits per month
- 100M+ total registered users
- Top 50 of the most visited sites worldwide





- Own low latency CDN
- 10K+ incoming simultaneous streams
- 1M+ outgoing simultaneous streams
- Try to keep minimal possible latency





🗙 NETINT













5



















Transcoding pipeline

Our approach:

• using only bare metal servers to achieve maximum network and resource utilization











9

Transcoding pipeline

Our approach:

- using only bare metal servers to achieve maximum network and resource utilization
- using high performance system profile for our servers











Transcoding pipeline

Our approach:

- using only bare metal servers to achieve maximum network and resource utilization
- using high performance system profile for our servers
- reducing resource wastage by implementing of all network parts (e.g. protocols for receiving and publishing streams) by ourselves













Transcoding pipeline

Our approach:

- using only bare metal servers to achieve maximum network and resource utilization
- using high performance system profile for our servers
- reducing resource wastage by implementing of all network parts (e.g. protocols for receiving and publishing streams) by ourselves
- using libav (which is part of ffmpeg) as a framework for transcoding inside our transcoder servers









📉 NETINT

(12)

Transcoding pipeline

Transcoder







13









Base	Dell PowerEdge R940xa	
Processor	(Intel Xeon Gold 6248 2.50Ghz 20 cores) x 4	
	24x16GB DDR4	
Memory	3200MHz ECC Reg	
Storage	2x480GB SSD	
	2x10G SFP+	
Network	(X710DA2)	



14







Transcoder	x264
Stream density 720p30 without source transcode	60-70













Problems:

• Number of full hd streams increased dramatically











Problems:

- Number of full hd streams increased dramatically
- Streamers started to use B-frames. One of our view protocols - WebRTC doesn't support B-frames







18

Software transcoder

Problems:

- Number of full hd streams increased dramatically
- Streamers started to use B-frames. One of our view protocols - WebRTC doesn't support B-frames
- Streamers started to increase GOP interval up to 10 second









Problems:

- Number of full hd streams increased dramatically
- Streamers started to use B-frames. One of our view protocols - WebRTC doesn't support B-frames
- Streamers started to increase GOP interval up to 10 second

















Transcoder	x264	
Stream density per server 1080p30	20	
Power consumption W	650	
Server cost \$	20`000	
Stream cost \$	1`000	
Production transcoding cost \$ (10K incoming streams)	10`000`000	













Transcoder	x264	
Stream density per server 1080p30	20	
Power consumption W	650	
Server cost \$	20`000	
Stream cost \$	1`000	
Production transcoding cost \$ (10K incoming streams)	10`000`000	





Production transcoding cost

Cs - server cost (with transcoding units, e.g. GPU cards) Ns - total number of streams NTS - capacity of the server *Round* - round up to the integer CTP - production transcoding cost

CTP = Round(Ns / NTs) * Cs







🗙 NETINT





Card	Stream density per card (1080p30)	Power consumption W	Card cost \$
Nvidia Tesla V100	20	250	8`000



23

- decoder/encoder queues must be around zero
- the cost of cards was valid at the time of tests (2020 year)







Card	Stream density per card (1080p30)	Power consumption W	Card cost \$
Nvidia Tesla V100	20	250	8`000
Nvidia Tesla P100	19	250	7`000



24

- decoder/encoder queues must be around zero
- the cost of cards was valid at the time of tests (2020 year)







Card	Stream density per card (1080p30)	Power consumption W	Card cost \$
Nvidia Tesla V100	20	250	8`000
Nvidia Tesla P100	19	250	7`000
Nvidia Tesla T4	16	70	3`000



- decoder/encoder queues must be around zero
- the cost of cards was valid at the time of tests (2020 year)





Transcoder	NVIDIA T4
Stream density per server 1080p30	96
Power consumption W	1`070
Server cost \$	38`000
Stream cost \$	395
Production transcoding cost \$ (10K incoming streams)	3`990`000



26





Transcoder	NVIDIA T4 (corrected to current market cost)	
Stream density per server 1080p30	96	
Power consumption W	1`070	
Server cost \$	32`000 (~2`000 per card)	
Stream cost \$	333	
Production transcoding cost \$ (10K incoming streams)	3`360`000	



27





ASIC transcoder















Card	T432	
Stream density per card (1080p30)	25	
Power consumption W	27	
Cost \$	2`000	



- decoder/encoder queues must be around zero
- the cost of cards was valid at the time of tests (2021 year)









Transcoder	T432	
Stream density per server 1080p30	150	
Power consumption W	812	
Server cost \$	32`000	
Stream cost \$	213	
Production transcoding cost \$ (10K incoming streams)	2`144`000	



30







Transcoder	T432 (corrected to current market cost)	
Stream density per server 1080p30	150	
Power consumption W	812	
Server cost \$	29`600 (1`600 per card)	
Stream cost \$	197	
Production transcoding cost \$ (10K incoming streams)	1`983`200	



31





Game changer?

Disadvantages:

- No hardware scaler. Need to scale in software.



32







Game changer?

Disadvantages:

- No hardware scaler. Need to scale in software.
- Can't reach real time (60 fps) for h264 4K60 fps.



33







Game changer?

Disadvantages:

- No hardware scaler.
 Need to scale in software.
- Can't reach real time (60 fps) for h264 4K60 fps

Even with software scaler on Dell R940 CPUs load is 40%. Cards limit is a bottleneck, NOT a CPU. Good solution if you have already had a park of servers with empty PCIe slots.





📉 NETINT



Quadra T2 transcoder



With hardware scaler inside!

Card	Quadra T2	
Stream density per card (1080p30)	45	
Power consumption W	40	
Cost \$	3`000	



• decoder/encoder queues must be around zero



Quadra T2 transcoder





Transcoder	Quadra T2	
Stream density per server 1080p30	270	
Power consumption W	890	
Server cost \$	38`000	
Stream cost \$	141	
Production transcoding cost \$ (10K incoming streams)	1`444`000	



36





Results: streams density

Туре	Card stream density	Server stream density (1080p30)	Number of servers
Software	0	20	500
NVIDIA T4	16	96	105
T432	25	150	67
Quadra T2	45	270	38

















Results: power consumption

Туре	Card power consumption W	Server power consumption W	Production power consumption W
Software	0	650	325`000
NVIDIA T4	70	1070	112`350
T432	27	812	54`404
Quadra T2	40	890	33`820















Results: cost

Туре	Card cost	Server cost \$	Stream cost \$	Production transcoding cost \$
Software	0	20`000	1`000	10`000`000
NVIDIA T4	2`000	32`000	333	3`360`000
T432	1`600	29`600	197	1`983`200
Quadra T2	3`000	38`000	141	1`444`000















(39)