

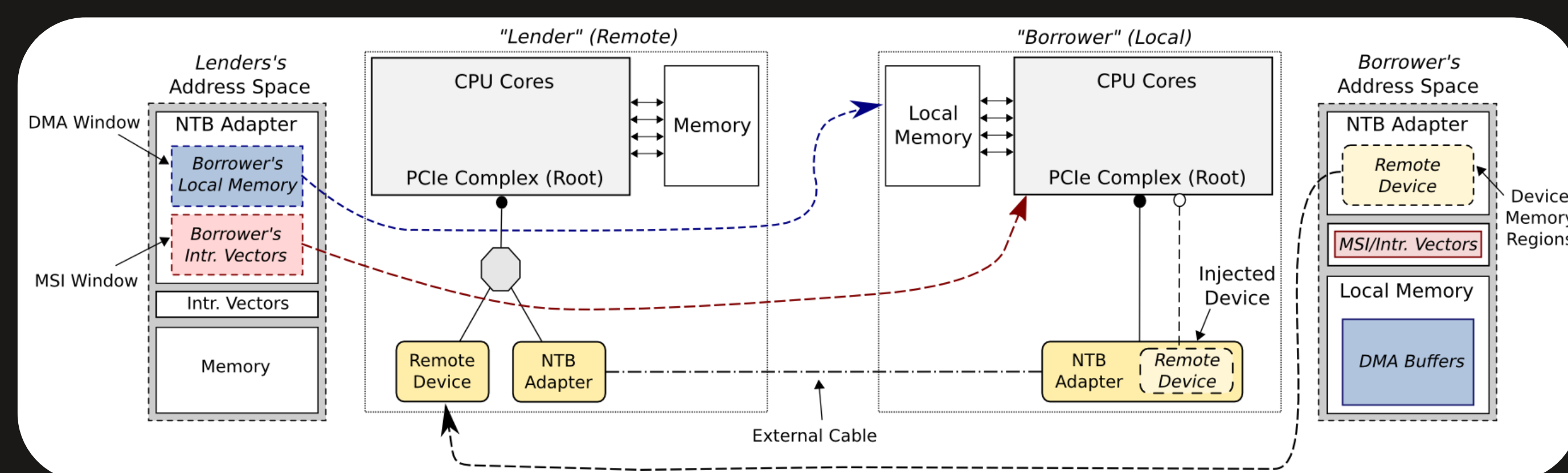
Efficient Processing of Medical Videos in a Multi-auditory Environment Using Gpu Lending

GPU LENDING OVER PCI-E

Transparent, low-latency, cross-machine.

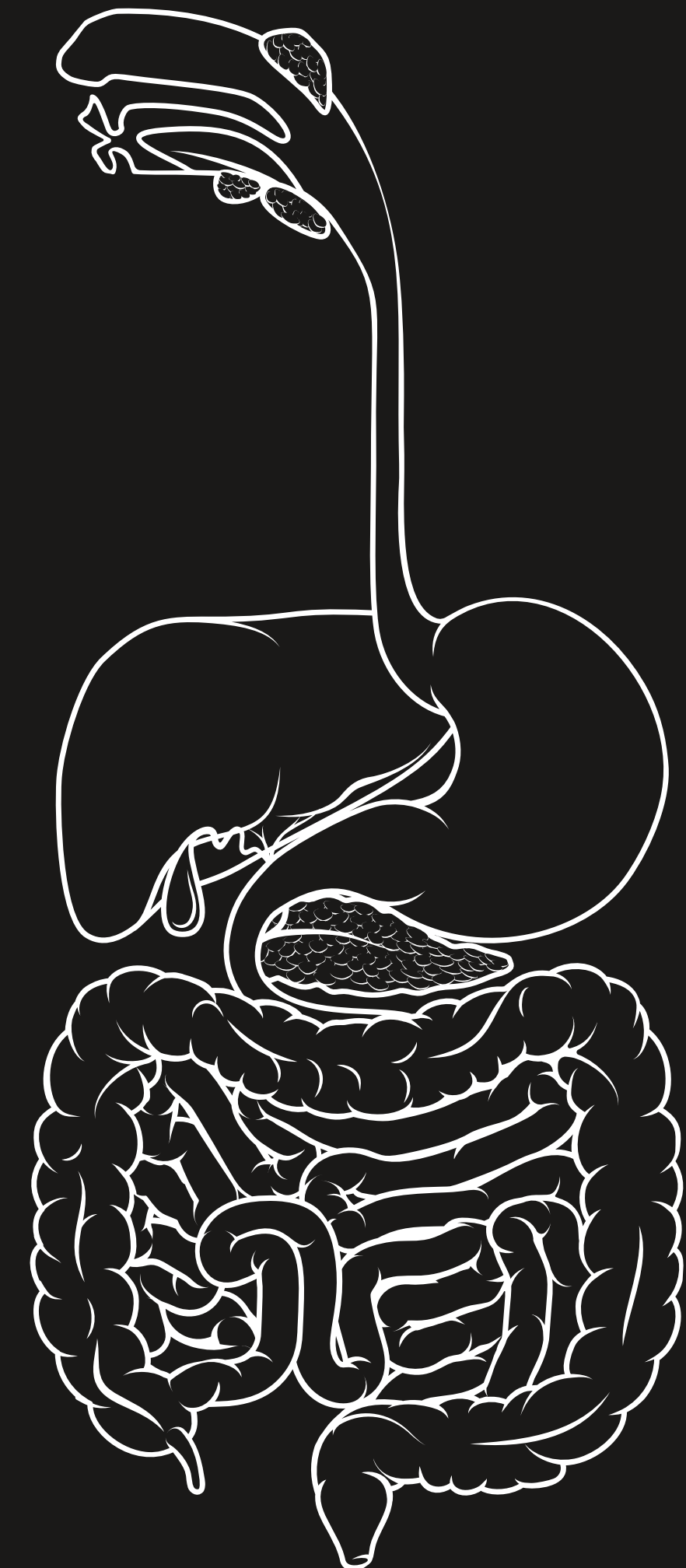
WE PRESENT DEVICE LENDING EXPERIMENTS RUNNING ON A SYSTEM DETECTING POLYPS IN THE GASTROINTESTINAL TRACT.

Computer-aided diagnostic systems (CADs) are designed to assist physicians in diagnosing patients. These systems are often based on machine learning and require GPU processing power. Such hardware is too noisy and requires space, so storing it in dedicated rooms is better. To preserve performance, we present experiments using Device Lending. Device Lending is implemented by Dolphin Interconnect Solutions PCI-e non-transparent bridges software and all devices connected to the network are considered part of one common resource pool.



WE ADDRESS TWO MAIN CHALLENGES:

- We show that real-time support is possible using this technology.
- We demonstrate the possibility of having one mainframe that lends the devices to different computers based on computational demands.



NO NEED FOR APPLICATION-SPECIFIC DISTRIBUTION MECHANISMS!

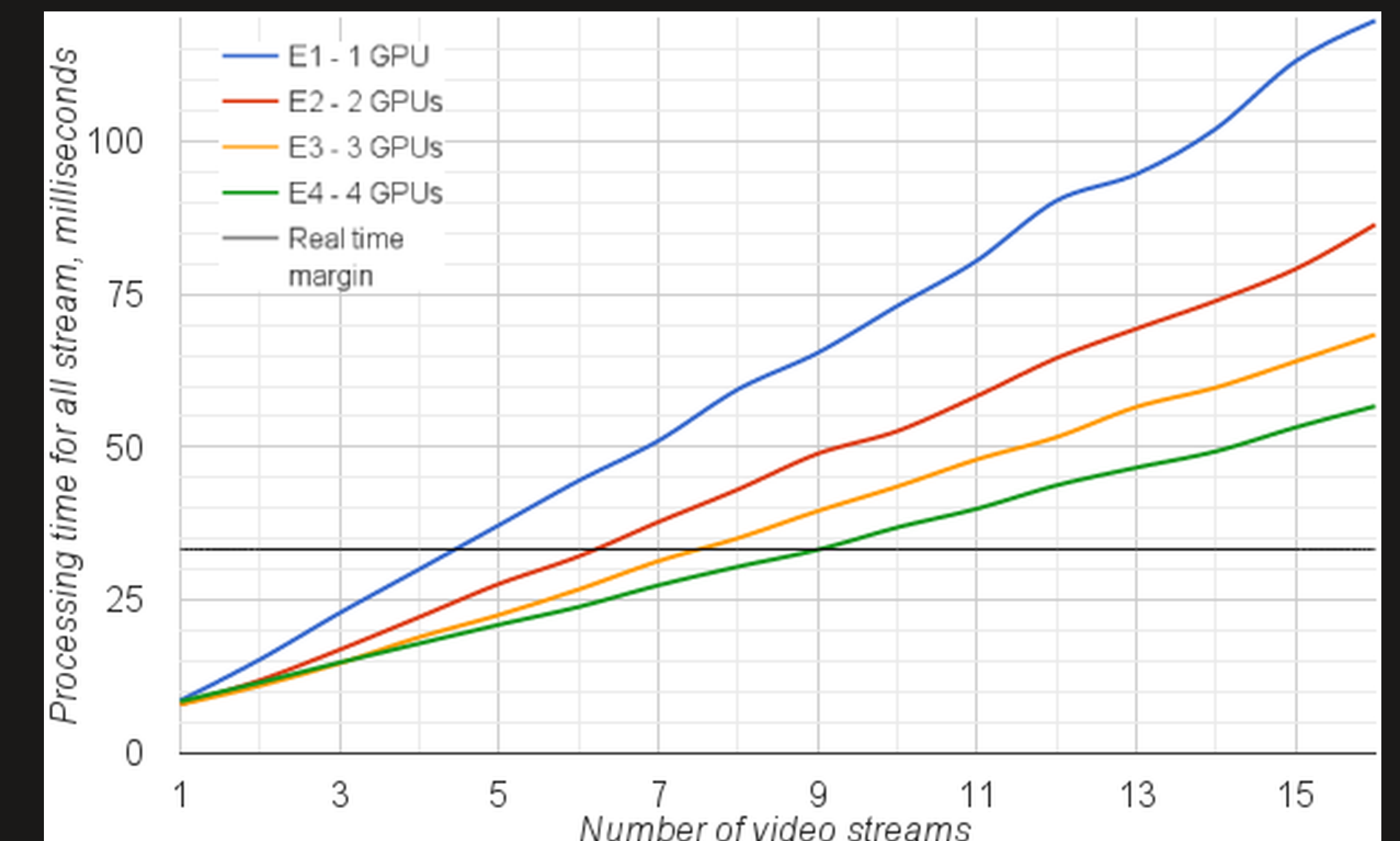
In the test, we have two computers, one with host code and the other lending a GPU. Experiment E1, E2, and E3 use local GPUs. In E4, we borrowed one GPU from the second computer in addition to three local GPUs. We performed polyp classification and real-time feedback on the video for up to 16 parallel video streams. We measured the performance from capturing the video up to showing the output on the screen.

WE SEE FROM THE EXPERIMENTS:

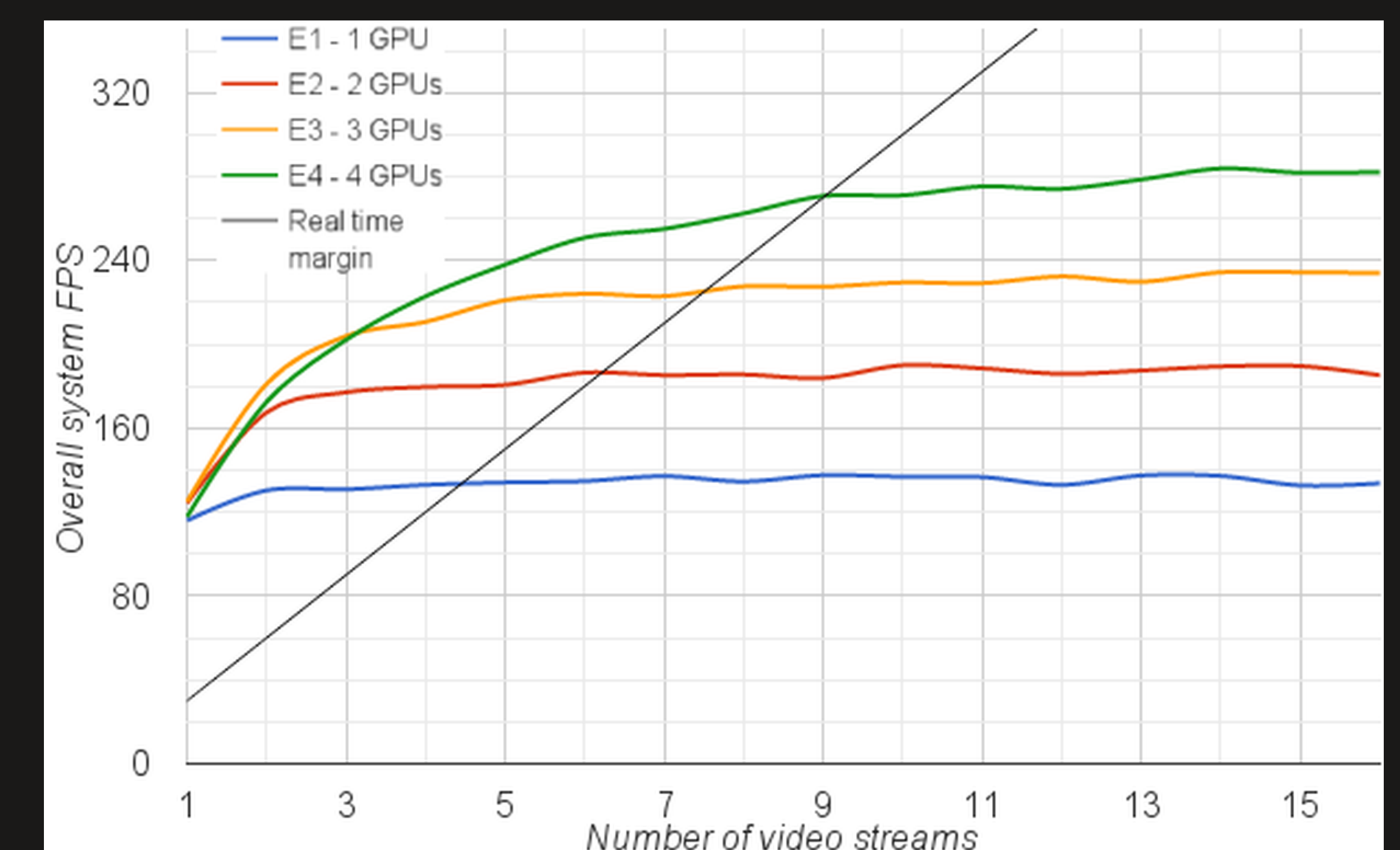
- Device Lending does not make sense for small amounts of data, but if the data to process is large it can give a large performance boost.
- Device Lending makes sense in our multi-auditory scenario.



With our initial experiments, we showed that for large-scale data in a multi-auditory environment Device Lending improved the performance. Aside from that, Device Lending can help in using resources and GPU power more efficiently. One challenge is that the planning of when a GPU is used in which room can create some overhead, but the positive aspects regarding power and resource efficiency are surpassing this. For future work, we would like to test our system with real physicians under clinical conditions. It would also be interesting to extend the idea to other scenarios such as cinemas and similar venues.



Frame processing time for several full HD streams in parallel.



Overall system performance for multiple full HD steams in parallel.

