# ON THE PLATFORN FOR OPTIMAL OVERLAY CHANNEL SELECTION IN  NETWORK POWERED by COMPUTING INVIRONMENT".

*R. Smelyanskiy, E. Stepanov*

*Moscow State University*

# Content

- **New Age for Computational Infrastructure**

- **Network Powered by Computing : concept**
- **InOpSys platform: Intelligent automatic network transport**
  **Optimization System**

- **Problem statement and  Modeling**

- **Experimental results**

# turing lecture

4.06.2018

**Innovations like domain-specific hardware, enhanced security, open instruction sets, and agile chip development will lead the way.**

BY JOHN L. HENNESSY AND DAVID A. PATTERSON

# A New Golden Age for Computer Architecture

# New Golden Age of Computational Infrastructure
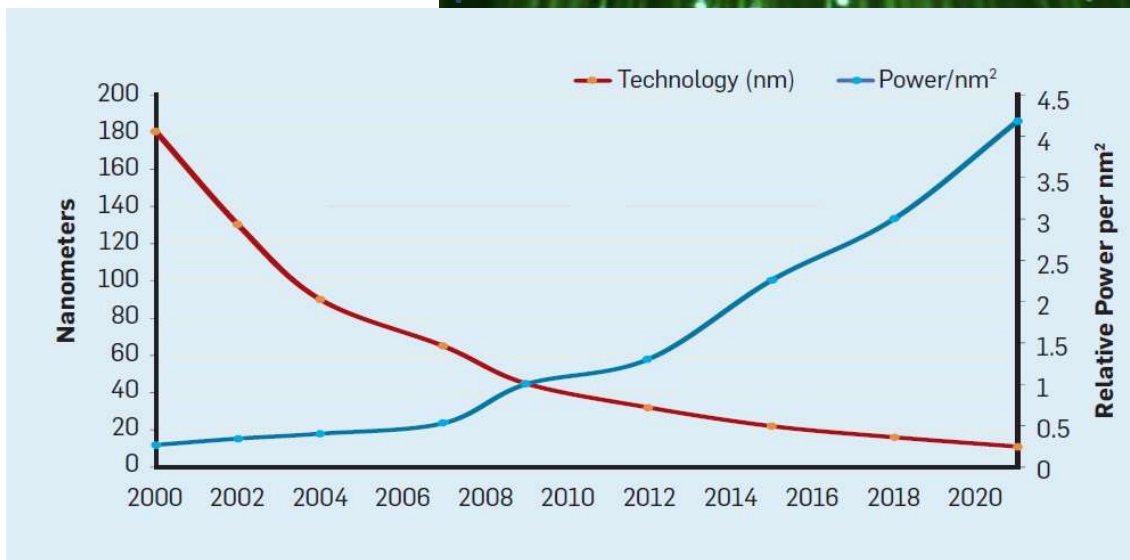
In 2022 the energy consumption of DC over the World was 200 GW x 24 hours x 365 days = 1 752 $10^{12}$ Wtph
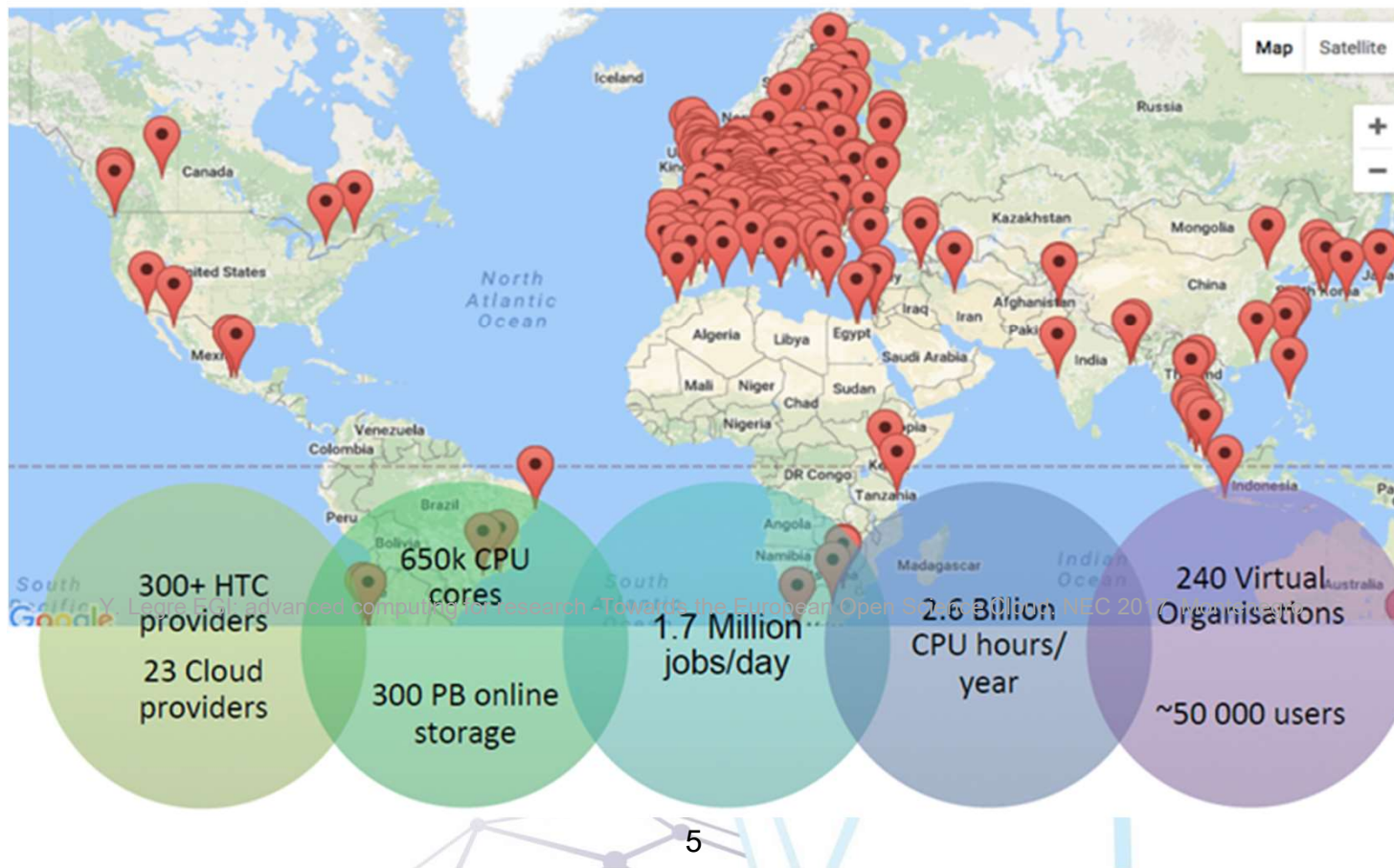
- 70-s - mainframe computer center with terminal



**Application Requirements + Hardware Capabilities + Software Engineering**

# EGI Federated Infrastructure



300+ HTC providers
23 Cloud providers

650k CPU cores
300 PB online storage

1.7 Million jobs/day

2.6 Billion CPU hours/year

240 Virtual Organisations
~50 000 users

# Applications suite of features

- **Distributed**
- **Real-Time mode**
- **Elasticity (SLA)**
- **Cross-platform**
- **Self-sufficiency**
- **Interaction and Synchronization**
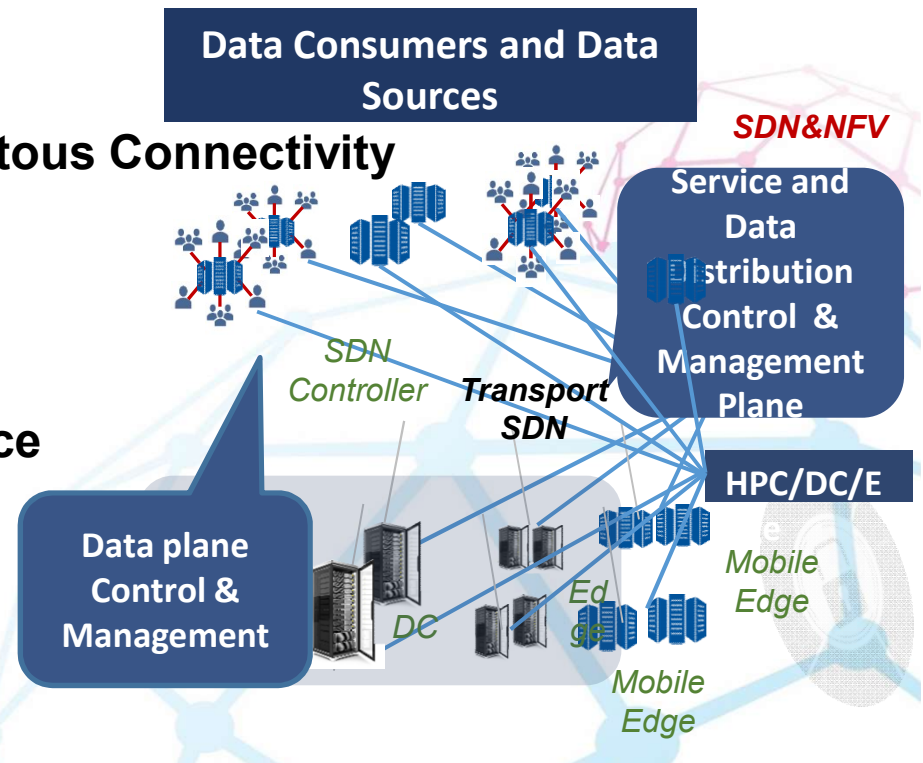- **Maintainability**

# Computational Infrastructure Requirements

- **Distributed Computing power and ubiquitous Connectivity**

- **Deterministic communication QoS**

- **Computing & Network Power Awareness**

- **Virtualization, Scalability, Serverless**

- **Availability, Reliability and Fault Tolerance**

- **Efficiency and Fairness**

- **Security**

**Data Consumers and Data Sources**

*SDN&NFV*

Service and Data Distribution Control & Management Plane

*SDN Controller*

*Transport SDN*

HPC/DC/E

Data plane Control & Management

*DC*

*Edge*
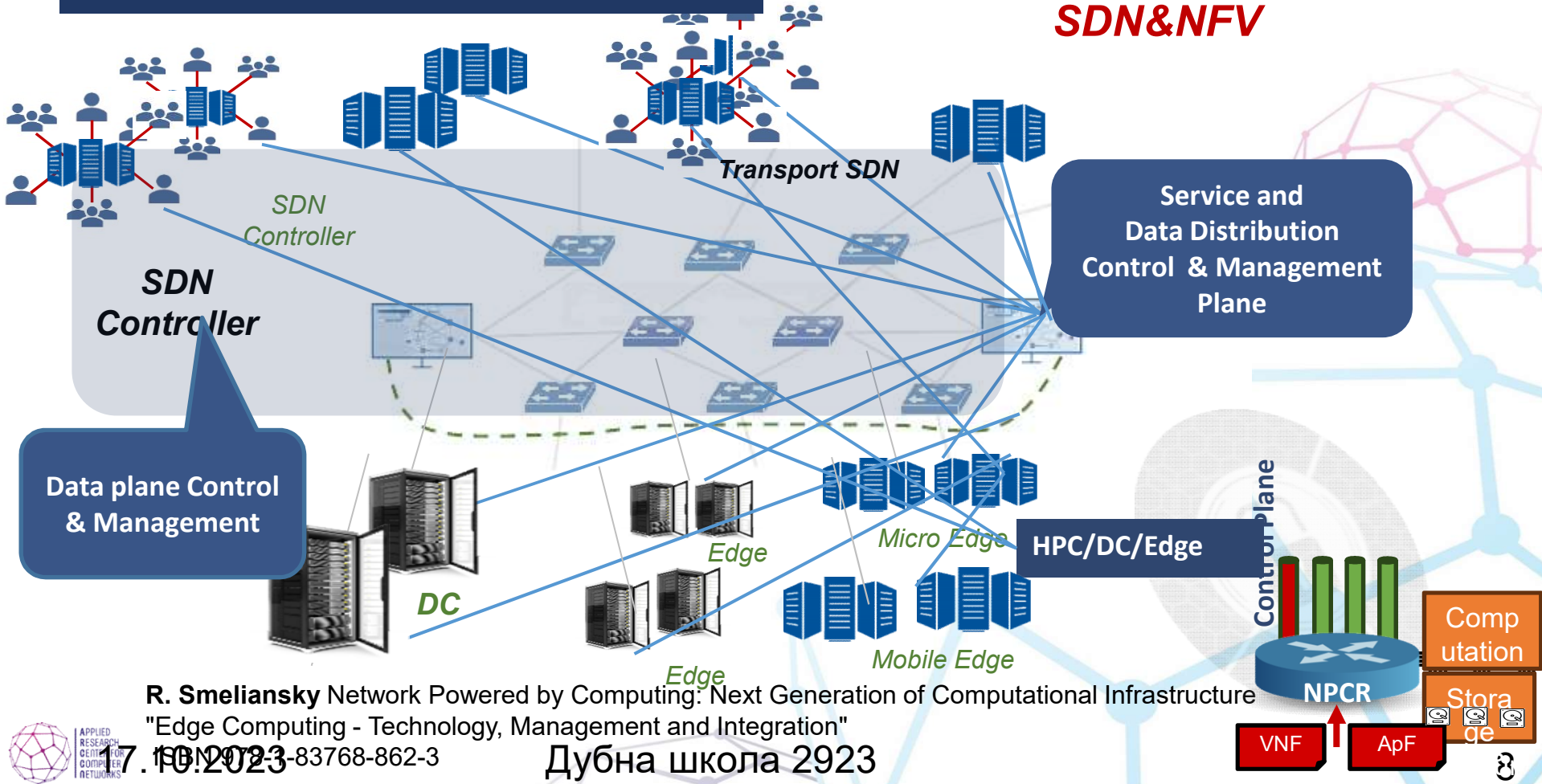
*Mobile Edge*

*Mobile Edge*

- **The scaling range of the network service is huge and in real time, which put high demands on the algorithm time complexity.**
- **Only sub-optimal solutions are available using methods based on machine learning**

**Data Consumers and Data Sources**

**SDN&NFV**

Transport SDN

SDN Controller

**SDN Controller**

Service and Data Distribution Control & Management Plane

**Data plane Control & Management**

*Edge*

*Micro Edge*

**HPC/DC/Edge**

Control Plane

*DC*

*Edge*

*Mobile Edge*

Computation

Storage

NPCR

VNF    ApF

**R. Smeliansky** Network Powered by Computing: Next Generation of Computational Infrastructure "Edge Computing - Technology, Management and Integration" ISBN 978-3-83768-862-3
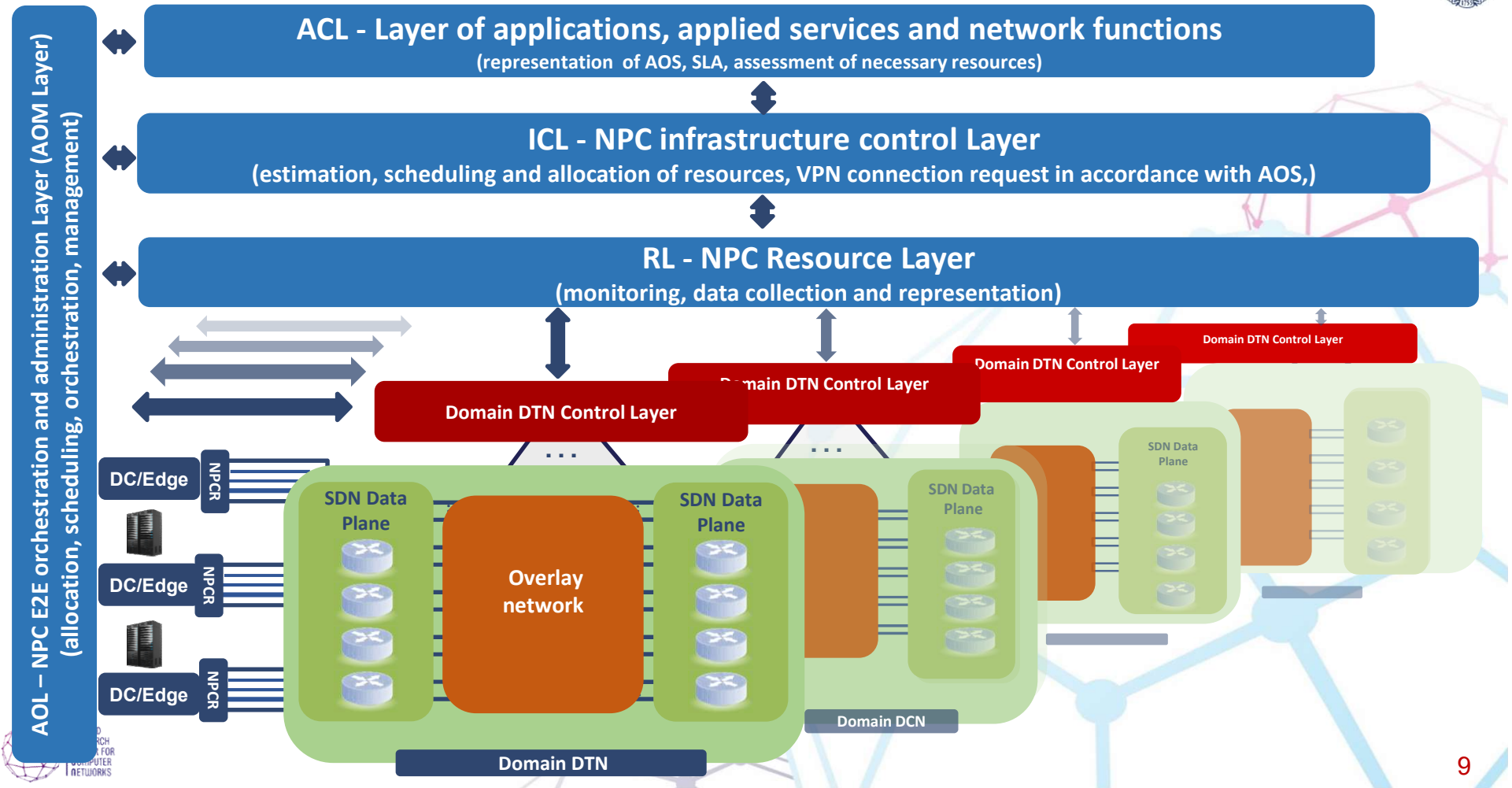
17.10.2023

Дубна школа 2923

8

# NPC Functional Architecture



**AOL – NPC E2E orchestration and administration Layer (AOM Layer)**
**(allocation, scheduling, orchestration, management)**

**ACL - Layer of applications, applied services and network functions**
(representation of AOS, SLA, assessment of necessary resources)

**ICL - NPC infrastructure control Layer**
(estimation, scheduling and allocation of resources, VPN connection request in accordance with AOS,)

**RL - NPC Resource Layer**
(monitoring, data collection and representation)

Domain DTN Control Layer

Domain DTN Control Layer

Domain DTN Control Layer

Domain DTN Control Layer

DC/Edge

NPCR

DC/Edge

NPCR

DC/Edge

NPCR

SDN Data Plane

Overlay network

SDN Data Plane

SDN Data Plane

SDN Data Plane

Domain DCN

Domain DTN

9

# NPC intra DTN Layer



**Control Plane Services**

Security OAM
(policy access, key management etc )

Connectivity control QoS control Monitoring

VPN channel distributed ledger

SDN controller

to the neighboring federate control plane

to the neighboring federate control plane

to the neighboring federate control plane

to the neighboring federate data plane

DC/Edge

NPCR

DC/Edge

NPCR

SDN Data Plane

SDN Data Plane

VPN Overlay Network

DTN Domain Data Plane

APPLIED RESEARCH CENTER FOR COMPUTER NETWORKS

# Problems Road Map

o  Implementation of ACL, ICL, NPC RL, NPC AOL functionality of Service and Data Distribution Management and Control Plane

o  **Optimal data traffic routing control**:
   o  **Selection of optimal overlay channel;**
   o  Data traffic balancing;

o  **Optimal allocation of application functions** (ApF)/ virtual network functions (VNF) across computational nodes (CN) of DP plane:
   o  ApF execution time estimation for certain CN;
   o  Selection of CN that optimal for execution of certain ApF/VNF

User request flows and ApS/Data Flows

Overlay channels/Data Plane

Control Plane

NPCR

Computational resources

Storage resources

VNF    ApF

# Intelligent automatic network transport Optimization System:            Problem description



Mpls 1

Conference app + FEC

A

Financial app + replication

Mpls 2

B

**Edge**   **NPCR**            **NPCR**   **Edge**

Public Internet

desktop app

SD-WAN DTN

**Flow SLAs:**

User Requirements: A→B
- Conference app:
  - Latency: 30ms, loss: 0.01, jitter: 10ms
- Desktop app:
  - Latency: 80ms, loss: 0.02, jitter: 30ms
- Financial app: high available
  - Latency: 40ms, No loss, no-jitter

**Techniques:**

Available Techniques:
- FEC:
  - Redundant pkt to fix packet loss
- TCP acceleration:
  - TCP proxy with new CC algorithm

**Links states:**

Three links between A and B:
- MPLS link 1:
  - Latency: 30ms, loss: 0.02, jitter: 10ms
- MPLS link 2:
  - Latency: 40ms, loss: 0.005, jitter: 5ms
- Internet link:
  - Latency: 70ms, loss: 0.01, jitter: 30ms

**User:**
- Input flow SLA requirement

**Prototype:**
- Monitor:
  - collet network states, links' KPI、CPU、memory etc.
- Evaluator:
  - if link can fit the flow, if not why?
- Policy generator:
  - Generate tech policies to improve the link

**Results:**
- Improve network to meet the SLA requirement

APPLIED RESEARCH CENTER FOR COMPUTER NETWORKS

# Optimization Problem: Statement

Each coming flow has SLA $\mathcal{A} = (B, D, J, L)$.
Each SD-WAN channel has current channel state :
$$S_r = (\underline{R_r}, \hat{R}_r, \overline{R}_r, N_r, l_r, \hat{j}_r, \overline{j}_r, \widehat{b}_r, \overline{b}_r, h_r)$$

subject to the SLA constraints:

$$\xi_b \ss(S_r, C) - B \geqslant 0$$

minimal admissible bandwidth,

$$D - \xi_d^1 \frac{(K + \Upsilon(C))p}{KB(S_r, C)} - \xi_d^2 \hat{R}_r \geqslant 0$$

maximal admissible time delay,

$$J - \xi_j \left(1 + B(h_r)^+\right) \overline{j}_r \geqslant 0$$

the maximal admissible jitter,

$$L - \mathcal{L}(S_r, C) \geqslant 0.$$

maximal admissible probability for packet loss.

The decision vector $C = (r, f, \gamma, c, \delta)$

$r \in \{1, \ldots, R\}$ – channel number
$f \in F$ – FEC algorithm
$\gamma \in \Gamma$ – FEC algorithm parameters
$c \in C$ – congestion control algorithm
$\delta \in \Delta$ – congestion control algorithm parameters
$\overline{l}_r$ - max. probability of retransmit.

$B$ – min admissible bandwidth
$D$ – max time delay
$J$ – max jitter
$L$ – max loss
$p$ – packet length
$\underline{R_r}, \hat{R}_r, \overline{R}_r$ - min, average, max RTT
$N_r$ - number of sent packets
$l_r$ - packet loss rate
$\hat{j}_r, \overline{j}_r$ - average, max jitter
$\widehat{b}_r, \overline{b}_r$ - average, max bandwidth
$h_r$ - current total load
$K$ – inf. pack. batch size
$\mathcal{E}_r = (T_r, C_r)$
$T_r, C_r$ - rent period and price
$M(C)$ – FEC batch size

$$\min_C \left(Q(S_r, \mathcal{E}_r, C) = \frac{C_r}{T_r} \cdot \frac{K + M(C)}{KB(S_r, C)} p \cdot (1 + \overline{l}_r)\right)$$

One packet transmission cost for the loop
«transmission – one possible retransmission» for channel $r$

13

# Packet Loss Codes

Suggested codes to use:

## 1-PR

$\oplus$

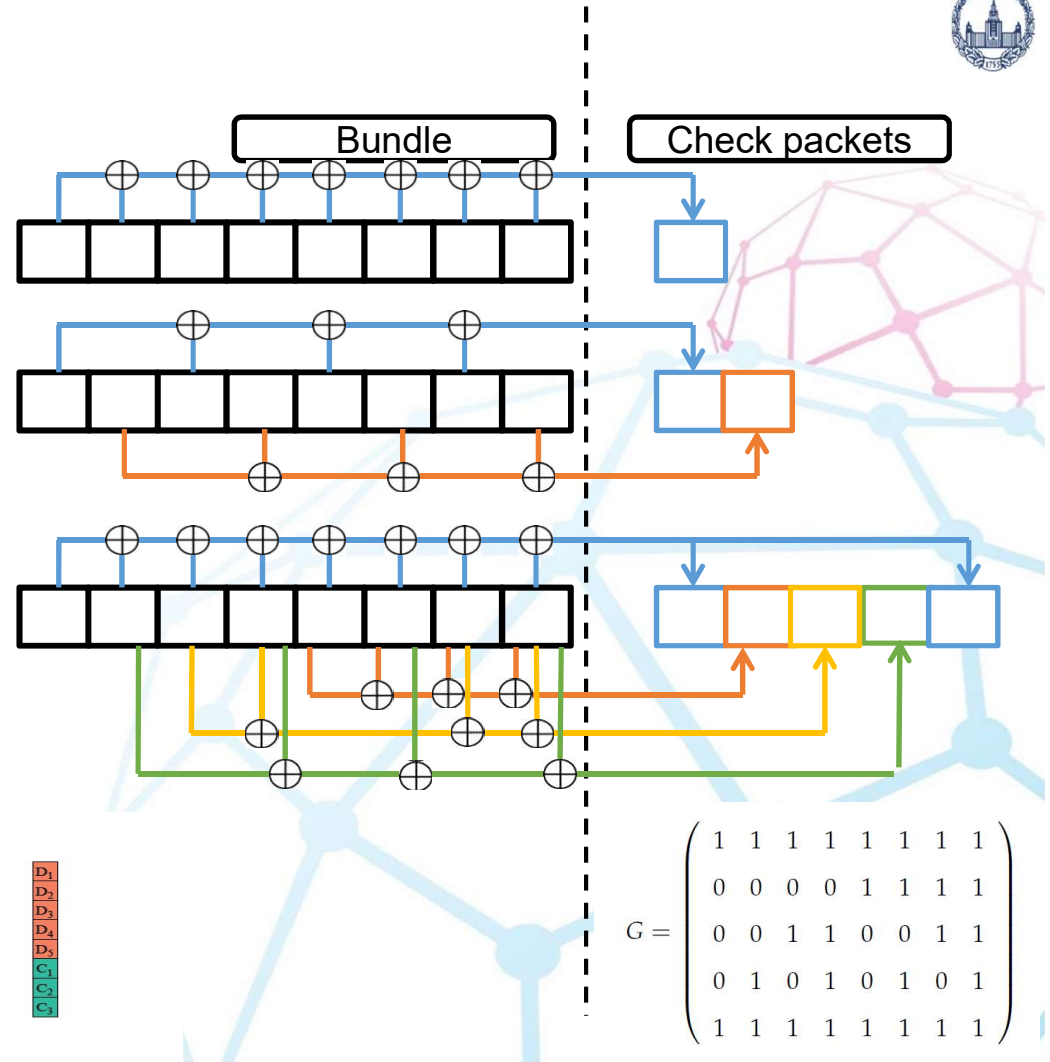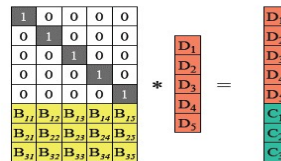- - XOR of all packets
- M=1
- Restores one packet

## 2-PR

- XOR for even and odd packets
- M=2
- Restores one packet in general and two packets if they are of different parity

## R-code

- Generating matrix of Reed-Muller code G with additional last line
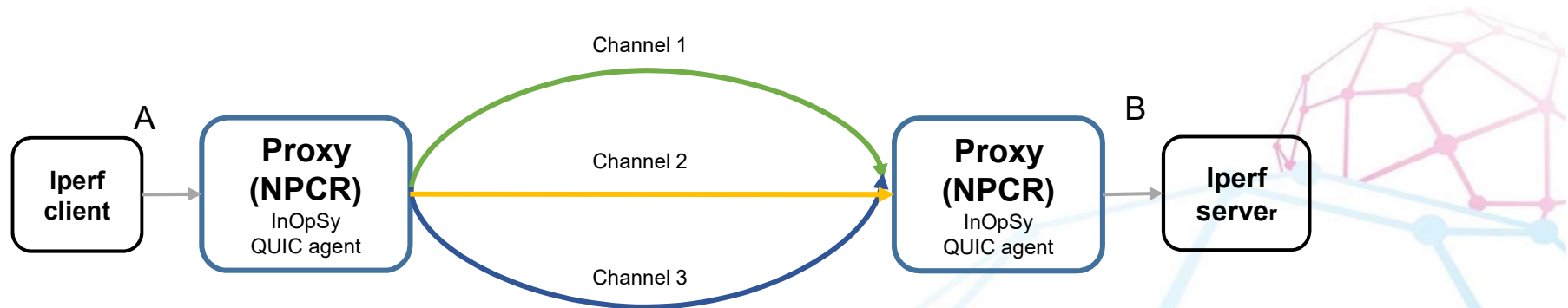- M= $\lceil \log_2 K \rceil + 2$
- Restores two packets

## Reed-Solomon Code (RS)



Bundle

Check packets

$$G = \begin{pmatrix} 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \\ 0 & 0 & 0 & 0 & 1 & 1 & 1 & 1 \\ 0 & 0 & 1 & 1 & 0 & 0 & 1 & 1 \\ 0 & 1 & 0 & 1 & 0 & 1 & 0 & 1 \\ 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \end{pmatrix}$$

# Testbed InOpSys platform



**Hardware:**

- Processor: 16 CPUs: Intel(R) Xeon(R) CPU E5-2667 v4 @ 3.20GHz
- Memory (RAM): 8 banks with DIMM DDR4 Synchronous 2133 MHz by 8GiB (64GiB total)
- OS: Ubuntu 18.04.6 LTS

**Channel QoS and SLA generation:**

- RTT: uniform distribution with parameters [10 ms, 100 ms].
- R: uniform distribution with parameters [25 Mbit/s, 50 Mbit/s].
- Loss: uniform distribution with parameters [0.00001 %, 5 %].
- J (jitter): uniform distribution with parameters [0.0 ms, 5 ms].

# Congestion control algorithm adjustment

The decision vector $\mathcal{C} = (r, f, \gamma, c, \delta)$, where

$r \in \{1, \ldots, R\}$ – channel number

$f \in F$ – FEC algorithm

$\gamma \in \Gamma$ – FEC algorithm parameters

$c \in C$ – congestion control algorithm

$\delta \in \Delta$ – congestion control algorithm parameters

Congestion control algorithm BBR parameters :

- BBRLossThresh = 2
- BBRBeta = 0.7
- BBRProbeRTTCwndGain = 0.5
- ProbeRTTDuration = 200 ms

## Bandwidth with BBR parameters



|  | Sender Speed (Kbit/s) | LOSS3 (%) | MinimalSpeed | MaximalSpeed | RTT ratio |
|---|---|---|---|---|---|
| Powell | 41846.8102 | 0.4067 | 40165.1476 | 43321.0271 | 1.3113 |
| Edges | 42208.2246 | 0.3934 | 40600.2356 | 43523.0272 | 1.3177 |
| Powell and Edges ratio | 1.0086 | 0.9674 | 1.0108 | 1.0047 | 1.0048 |

16

# ML modeling $\mathcal{B}(S_r, \mathcal{C})$ and $\mathcal{C}$ parameters

Find $\mathcal{B}(S_r, \mathcal{C})$- the bandwidth that can be reached on the SD-WAN channel with the current state Sr and chosen decision vector $\mathcal{C}$ – parameters of CC algorithm

$Find\ \mathcal{C} = (r, f, \gamma, c, \delta)\ , where$
$\delta \in \Delta$ – congestion control algorithm parameters: BBRLossTresh, BBRBeta, BBRProbeRTTCwndGain, ProbeRTTDuration



Dataset size vs model quality = deviation model vs experiment.

# ML models Ensemble

- $\alpha$ is changing in [0, 1]

- $\beta$ is changing in [0, 1 - alpha]

Then, our ensemble have following structure:

$$MODEL1\_PRED * \alpha + MODEL2\_PRED * \beta + MODEL3\_PRED * (1 - \alpha - \beta)$$

Changing alpha, beta parameter to see, how models impact in overall result.



Alpha = 0.85
Beta = 0.1

$ML \; \mathcal{B}(S_r, \mathcal{C})$ model accuracy ~ 0.988

# Experiments examples (bandwidth)

| Channels (BW-1 experiment) | | | |
|---|---|---|---|
| | Mean BW (Mbit/s) | Mean RTT (ms) | Mean Jitter (ms) | Loss (%) |
| №0 | 100 | 30 | 1 | 0.01 |
| №1 | 120 | 30 | 1 | 0.01 |
| №2 | 100 | 30 | 1 | 0.01 |

| SLA | | | | |
|---|---|---|---|---|
| BW (Mbit/s) | One-way Delay (ms) | Jitter (ms) | Loss (%) | Chosen Channel |
| 80 | 16 | 3 | 0.001 | №1 |

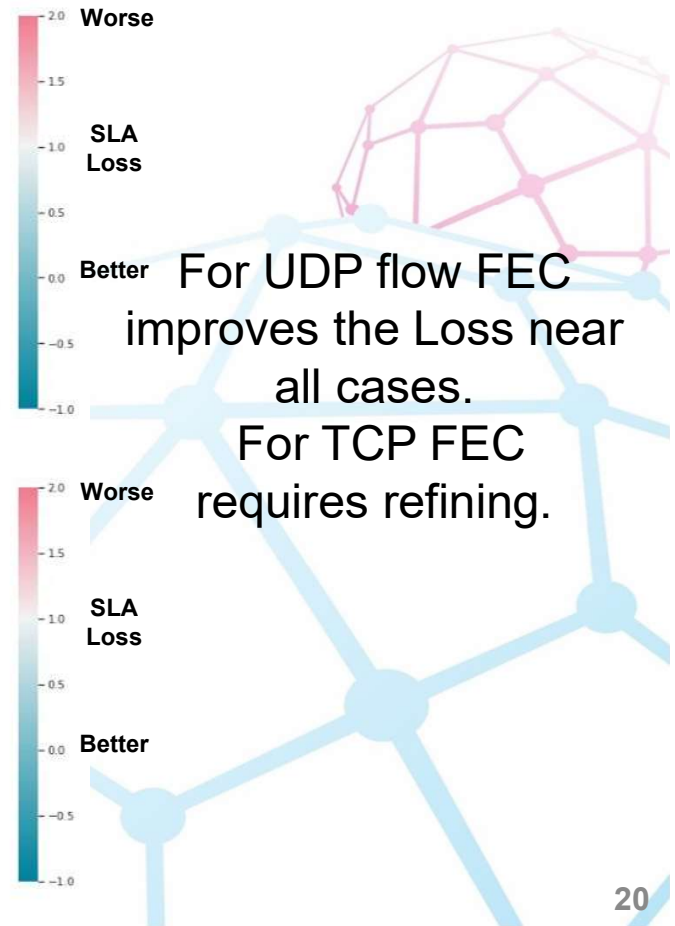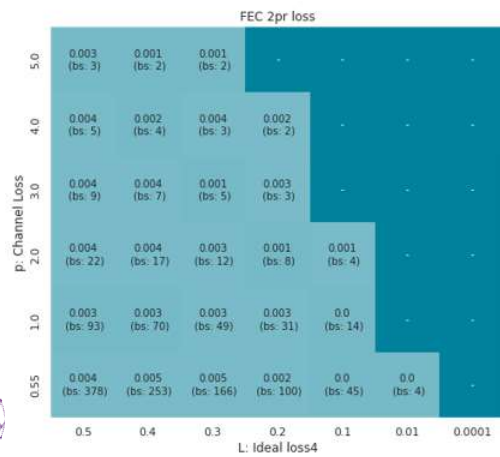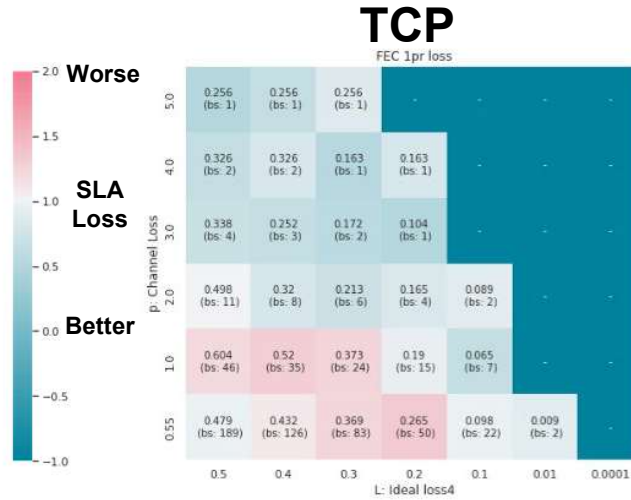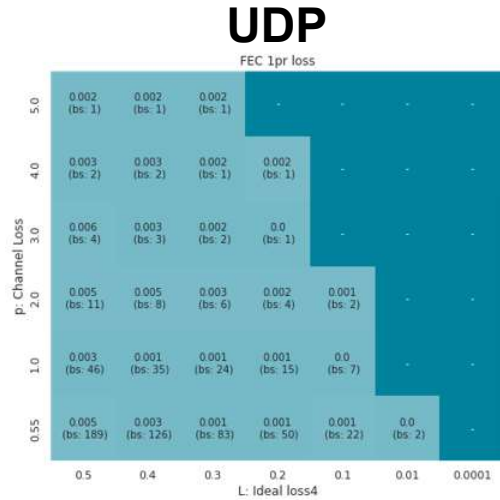BW-1: Three channels match required SLA, sometimes in excess

| Channels (BW-2 experiment) | | | |
|---|---|---|---|
| | Mean BW (Mbit/s) | Mean RTT (ms) | Mean Jitter (ms) | Loss (%) |
| №0 | 80 | 30 | 1 | 0.01 |
| №1 | 100 | 30 | 1 | 0.01 |
| №2 | 80 | 30 | 1 | 0.01 |

| SLA | | | | |
|---|---|---|---|---|
| BW (Mbit/s) | One-way Delay (ms) | Jitter (ms) | Loss (%) | Chosen Channel |
| 81 | 16 | 3 | 0.001 | №1 |

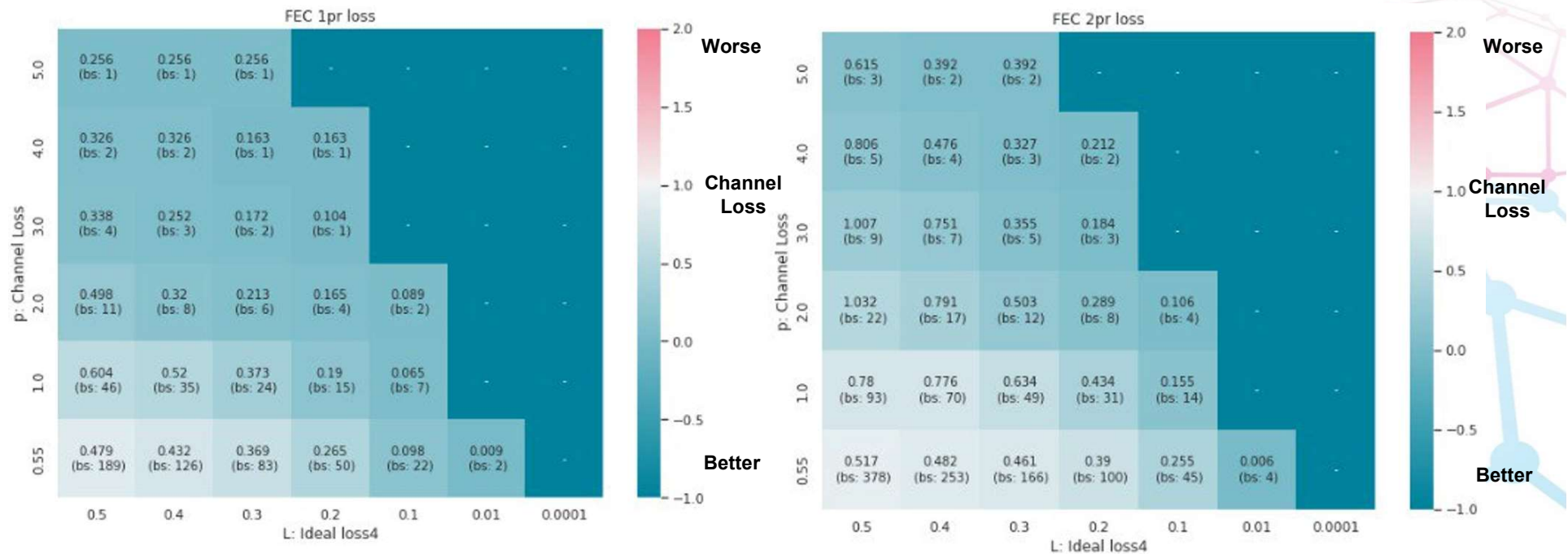BW-2: Second channel fits SLA, others are out of SLA because of low bandwidth.

# Loss evaluation: UDP vs TCP

**UDP**

**TCP**

For UDP flow FEC improves the Loss near all cases.
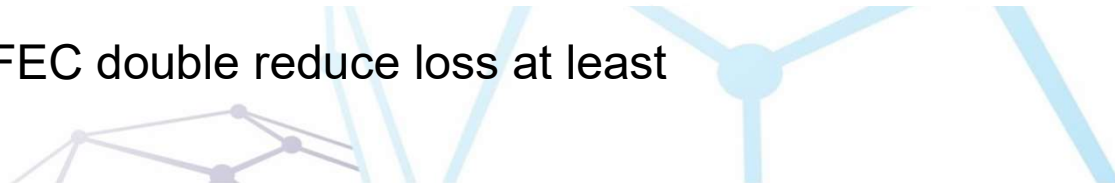For TCP FEC requires refining.

# Loss improving: TCP



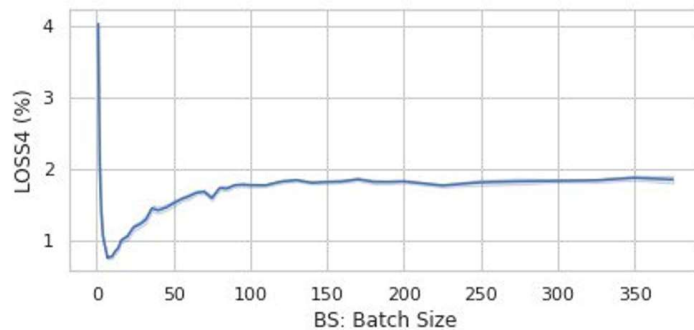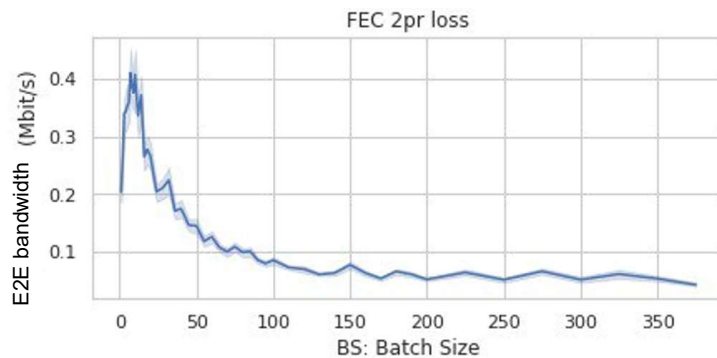For TCP  FEC double reduce loss at least

# Why the proposed FEC mechanism does not improve the loss parameter for TCP better then double?

## Congestion control window size is less than FEC batch size (zero packets) sometimes



FEC 2pr loss

There is an optimum batch size that gives the highest bandwidth and the lowest loss

It is promising to apply machine learning methods to find the optimal value

# SLA compliance with adjusted coefficients

Constraints:

$$\xi_b \mathcal{B}(\mathcal{S}_r, \mathcal{C}) - B \geqslant 0,$$

$$D - \xi_d^1 \frac{(K + M(\mathcal{C}))p}{KB(\mathcal{S}_r, \mathcal{C})} - \xi_d^2 \widehat{R}_r \geqslant 0,$$

$$J - \xi_j \left(1 + B(h_r)^+\right) \overline{j}_r \geqslant 0,$$

$$L - \mathcal{L}(\mathcal{S}_r, \mathcal{C}) \geqslant 0.$$

$$\xi_j = 1.36, \ \xi_d^1 = 2.53, \ \xi_d^2 = 1.27 \ \xi_b = \begin{cases} 0.31, Loss < 0.1\%, \\ 0.71, elsewise \end{cases}$$

| Algorithm | BW SLA | Loss SLA | RTT SLA | Jitter SLA |
|---|---|---|---|---|
| InOpSys | 96% | 97% | 90% | 100% |

| | $\dfrac{RTT_{simulation}}{2D}$ | $\dfrac{B_{simulation}}{B}$ |
|---|---|---|
| mean | 0,7911 | 1,3548 |
| 10% | 0,5471 | 2,0611 |
| 20% | 0,6037 | 1,3217 |
| 30% | 0,6518 | 1,2879 |
| 40% | 0,6977 | 1,2500 |
| 50% | 0,8043 | 1,2119 |
| 60% | 0,8711 | 1,1905 |
| 70% | 0,8983 | 1,1716 |
| 80% | 0,9335 | 1,1253 |
| 90% | 0,9991 | 1,081 |
| all | 1,4329 | 0,7653 |

worse

better

By adjustment the InOpSys platform meet SLA for more than 90% of the flows

# Comparative analysis

| Algorithm | Correct channel | RTT ratio | Loss ratio | QUIC speed ratio | End-to-End speed ratio |
|---|---|---|---|---|---|
| InOpSys | 100% | 0,99 | 0,62 | 1,49 | 1,08 |
| InOpSys (without FEC) | 100% | 1,43 | 0,75 | 0,89 | 0,73 |
| Random | 13,70% | 2,09 | 7,64 | 0,71 | 0,49 |
| Min RTT | 12,60% | 1,25 | 5,43 | 0,89 | 0,62 |
| Min Loss | 15% | 2,31 | 0,13 | 1,00 | 0,70 |
| vQoE | 14,60% | 1,26 | 3,65 | 0,91 | 0,62 |

| Algorithm | BW SLA | RTT SLA | Loss SLA | Total SLA |
|---|---|---|---|---|
| InOpSys | 96% | 90% | 97% | 90% |
| Random | 9,50% | 65% | 46% | 4,70% |
| MinRTT | 35,70% | 90,50% | 74,70% | 25,20% |
| MinLoss | 49,30% | 61,60% | 97,20% | 35,60% |
| vQoE | 42,20% | 90% | 76,60% | 33,30% |

better        worse

$$vQoE = vQoE'\left(\frac{BaselineLoss}{Loss}\right) + vQoE'\left(\frac{BaselineRTT}{RTT}\right)$$

InOpSy platform selects the channel with maximal injection speed and relatively small mean loss ratio.

# Overheads: "free cheese only in a mousetrap"

## Channel selection algorithm overhead

| Algorithm channel selection algorithm | Execution time (ms) | RSS (Kbyte) | VMS (Kbyte) | CPU (%) |
|---|---|---|---|---|
| InOpSys | 70,71 | 48,28 | 530,29 | 379,80 |
| Random | 0,02 | 47,88 | 402,78 | 98,70 |
| MinRTT | 0,02 | 47,88 | 402,78 | 99,60 |
| MinLoss | 0,03 | 47,89 | 402,76 | 99,50 |
| vQoE | 0,03 | 47,88 | 402,78 | 99,90 |

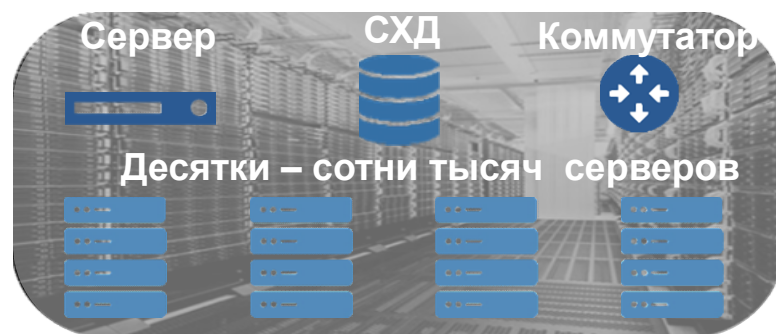Resident Set Size (RSS) – the memory volume available for process in RAM.

Virtual Memory Size (VMS) is the total volume memory available a process can access

## QUIC Agent operation overheads

| | Client | | | Server | | |
|---|---|---|---|---|---|---|
| | VMS (Kbyte) | RSS (Kbyte) | CPU (%) | VMS (Kbyte) | RSS (Kbyte) | CPU (%) |
| R-Scheme | 14776,00 | 6016,00 | 72,30 | 15094,00 | 8284,00 | 44,70 |
| Reed-Solomon | 15800,50 | 14900,00 | 44,70 | 16118,50 | 16796,00 | 43,10 |
| 1PR | 14776,00 | 6176,00 | 45,00 | 15094,00 | 8288,00 | 42,30 |
| 2PR | 14776,00 | 6100,00 | 47,80 | 15094,00 | 8236,00 | 45,00 |
| No FEC only batch | 4251,00 | 7668,00 | 40,70 | 4513,00 | 7608,00 | 40,50 |
| No FEC no batch | 4249,00 | 7484,00 | 46,50 | 4513,50 | 7736,00 | 44,80 |

better                                          worse

# Существующий подход к организации ЦОД



Облачное пространство

Виртуальная инфраструктура

Виртуальный сервис     Виртуальная СХД

Сервер     СХД     Коммутатор

Десятки – сотни тысяч серверов

Филиалы, Подразделения, Дочерние предприятия

А.В.

**Проблемы существующего подхода:**

- **Высокие требования к QoS каналов связи**, для обеспечения доступности сервиса;

- **Проблемы капитального строительства** централизованного ЦОД;

- **Проблема масштабирования**, связанная со строительством новых ЦОД и ручной синхронизацией их работы;

- **Неоптимальное использования доступных ресурсов**, из-за отсутствия централизованной системы управления и

APPLIED RESEARCH CENTER FOR COMPUTER NETWORKS

# Новый подход к организации ЦОД



Облачное
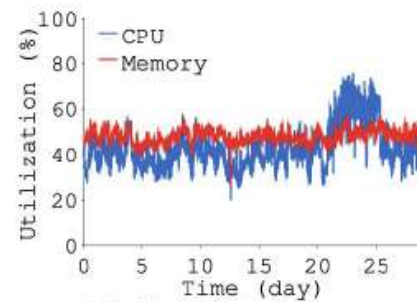пространство

Виртуальная
инфраструктура
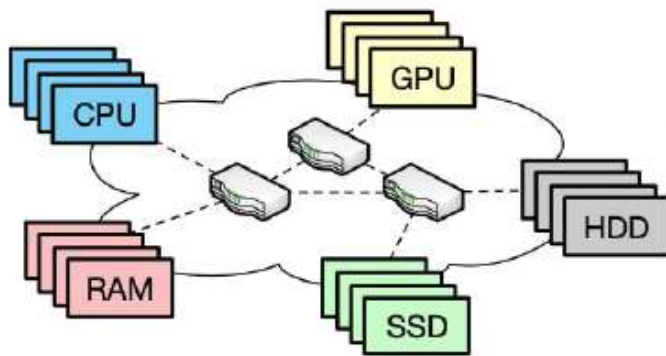
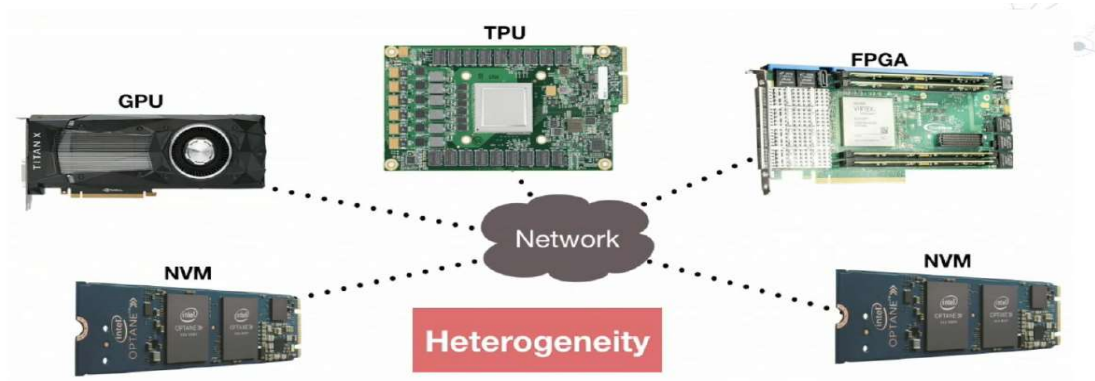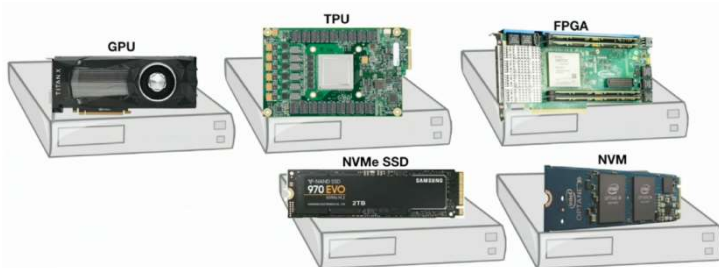Виртуальный сервис    Виртуальная СХД

Сеть микро—ЦОД'ов
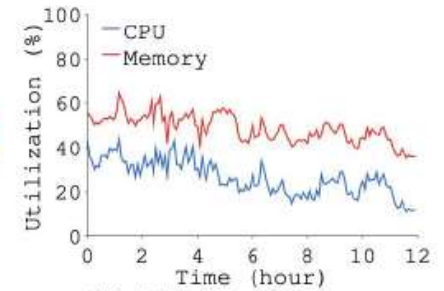
А.В.

**Преимущества нового подхода**

- **Снижение** требований к транспорту за счет близости экземпляра сервиса к конечному потребителю;
- **Снижение затрат** на организацию ЦОД за счет отсутствие необходимости строить централизованный ЦОД;
- **Простое масштабирование** за счет использования централизованной облачной платформы;
- **Повышение оперативности работы сети** за счет централизованной системы управления и оркестрирования и близости сервиса к клиенту.

27

# Дезагрегированная Архитектура
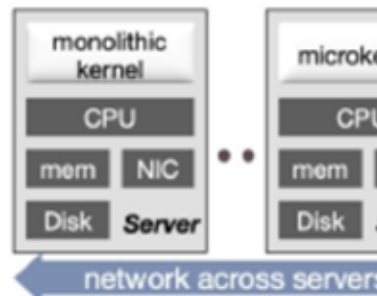


(a) Google Cluster

(b) Alibaba Cluster

**20% - 60% использования CPU и пам**

# Эволюция серверной архитектуры



**Монолитная**

monolithic kernel
CPU
mem | NIC
Disk | *Server*

microk...
CPU
mem
Disk

network across servers

40x B
640G
40nm

2010    2012    2014    2016    2017    2019    2021

Japan 2022 - 1,520,000 Gbps on Fiber-optic network
https://techxplore.com/news/2022-11-petabit-transmission-mode-fiber-standard.html

https://www.techpowerup.com/262237/broadcom-ships-25-6-tbps-network-switch-on-7-nm-chip#g2

# Conclusion

- **Growth of Application requirements are the big challenges for Computational Infrastructure management and control**

- **Network Powered by Computing Environment**
  - **Distributed Computing power and ubiquitous Connectivity**
  - **Deterministic communication QoS**
  - **Computing & Network Power Awareness**

- **InOpSys (Intellectual transport Optimization System) platform for NPC is presented that allows automatically:**

  - **select the best channel from the available ones by the metrics takes into account Bandwidth, Delay, Loss, Jitter**

  - **adjust the parameters of selected channel to meet the SLA requirements**

**National TE Data Sets for ML methods have to be developed**

**ML technics enable NPC environment to be predictable, secure, reliable, efficient, scalable.**

30

# THANKS

**Contacts: smel@cs.msu.su**