



计算机科学与工程学院
School of Computer Science and Engineering



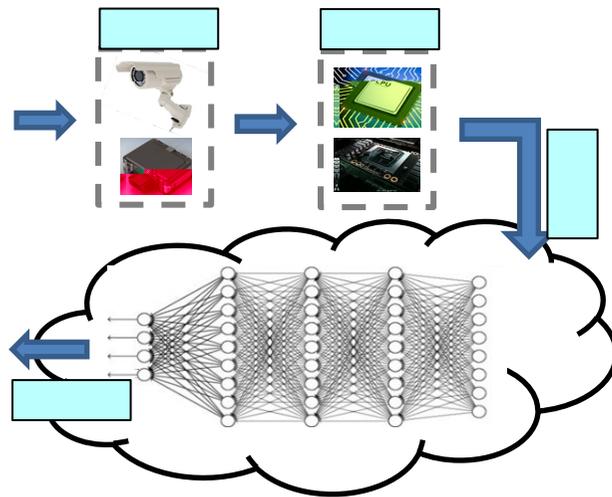
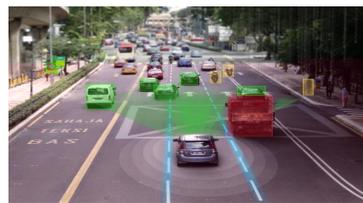
dshen@seu.edu.cn



□ ChatGPT

□

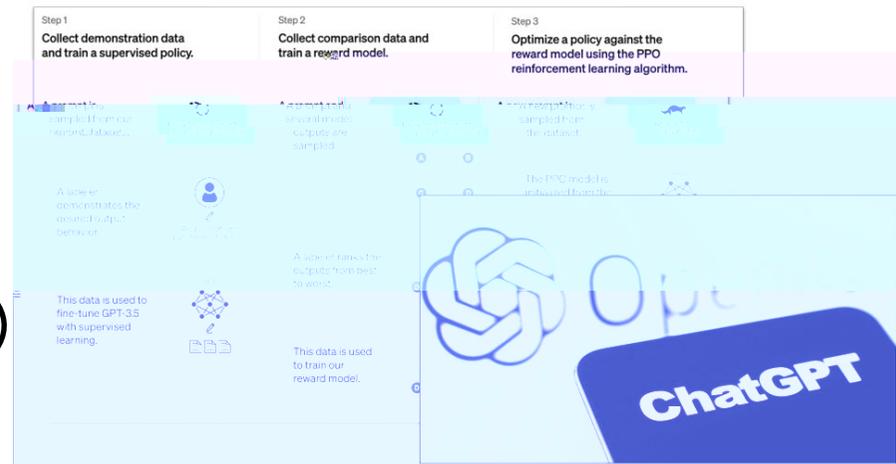
□



ms

[2]

Training Type	Cluster Configuration	Training Time
Fine-tune	A100-80G*8	One week
Reward Model	A100-90G*24	One week
Train ChatGPT	A100-80G*160	One week

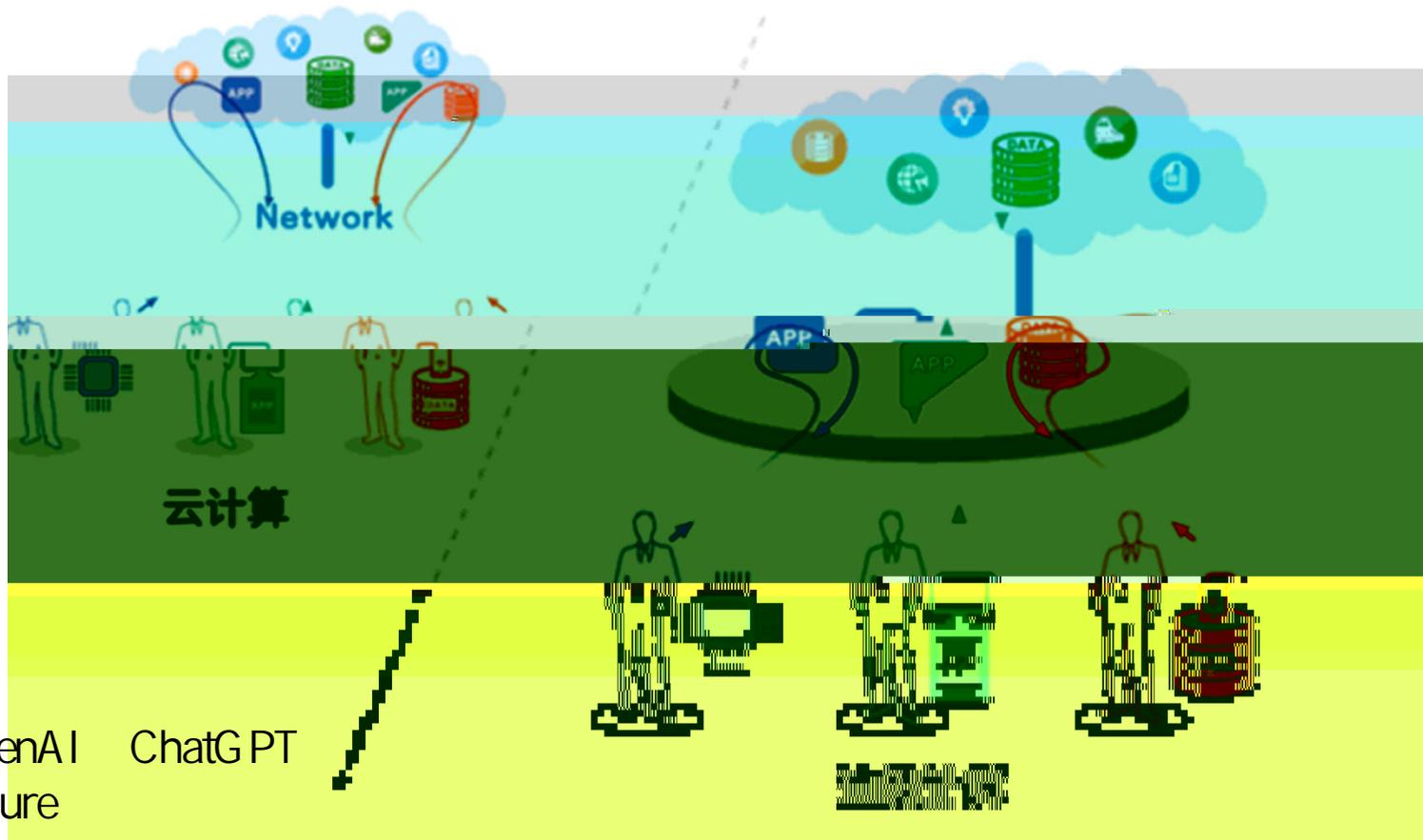


ChatGPT

[1]

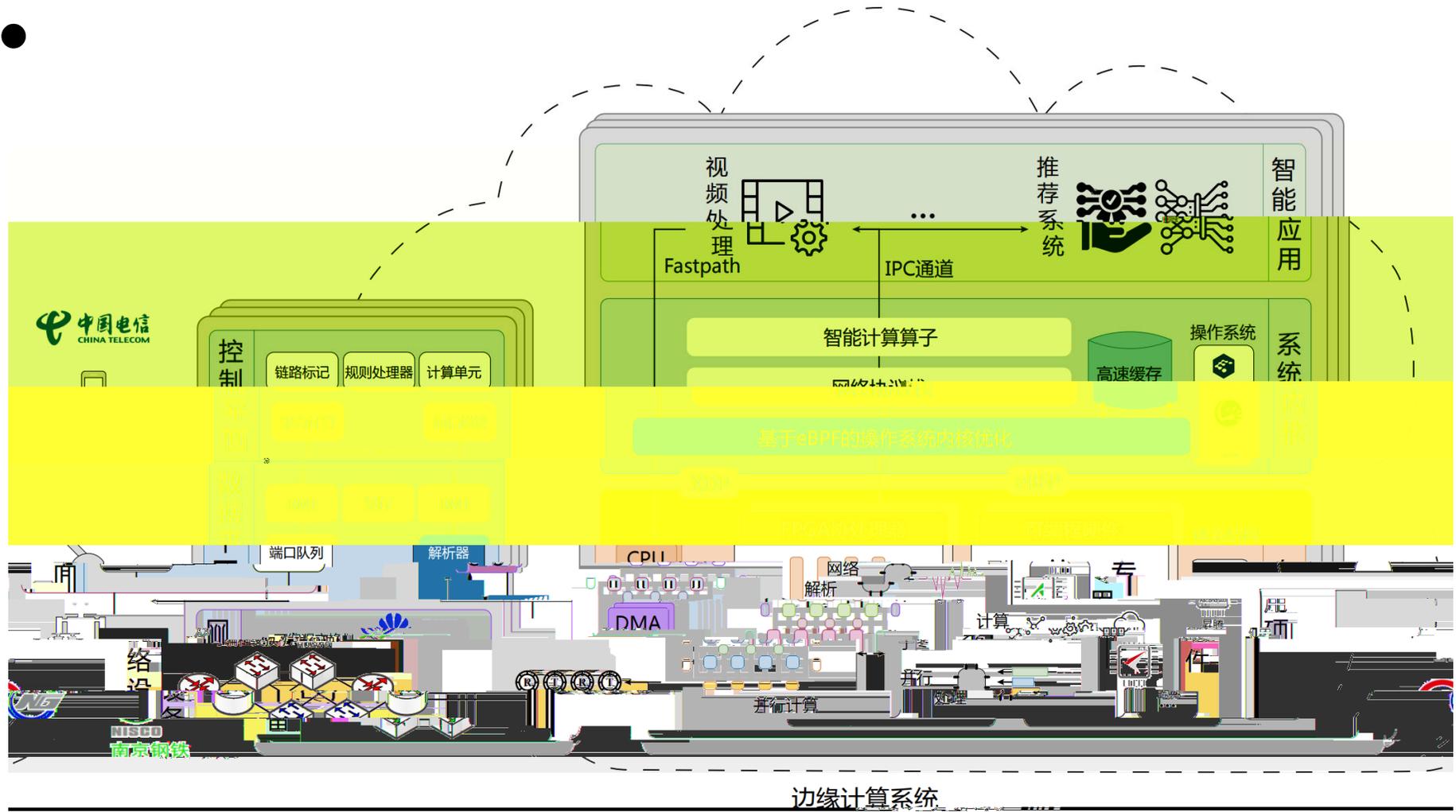
[1] Alibaba Cloud. Towards AI-Oriented High Computing Power Data Center Network. A PN et 2023.

[2] 2030 —"



Tips OpenAI ChatGPT
Azure

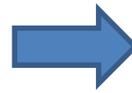
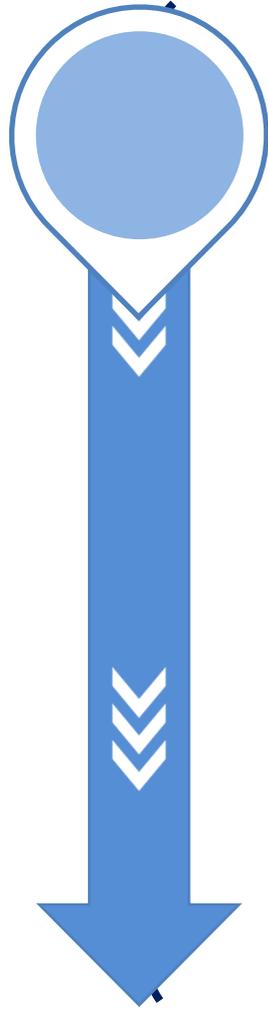
/





Linux eBPF

FPGA



Multi-exit DNN Inference Acceleration based on Multi-Dimensional Optimization for Edge Intelligence. IEEE Transactions on Mobile Computing, 2022 **CCF A**

Exploiting the Computational Path Diversity with In-network Computing for MEC. IEEE SECON 2022 **CCF B**

WAEVSR: Enabling Collaborative Live Video Super-Resolution in Wide-Area MEC Environment. IWQoS 2023 **(CCF B)**



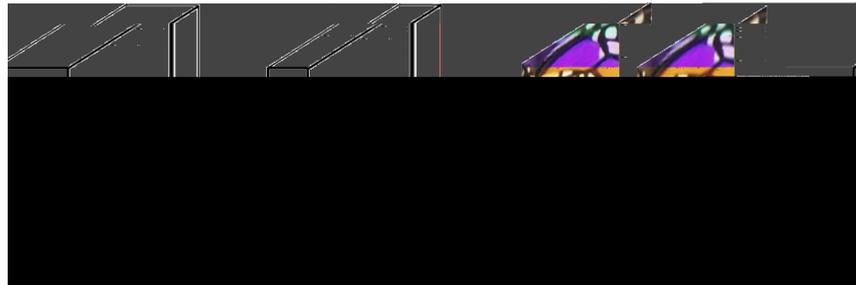
e.g. BCHW=7x 3x 480x 720



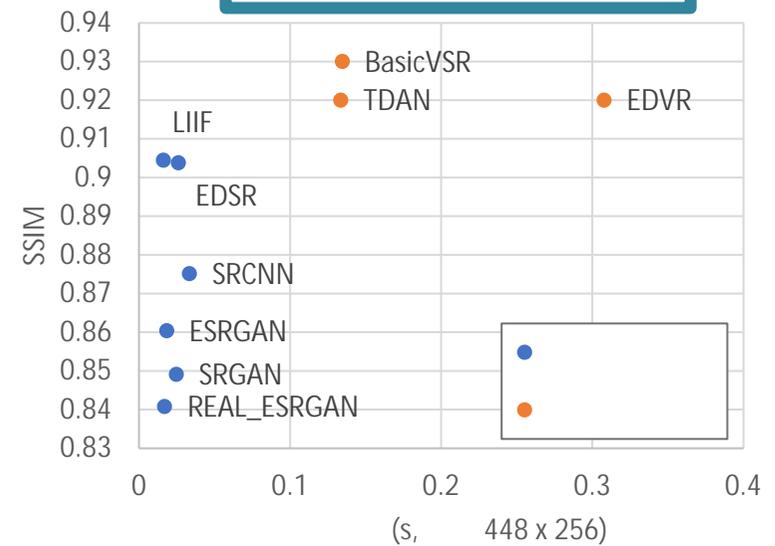
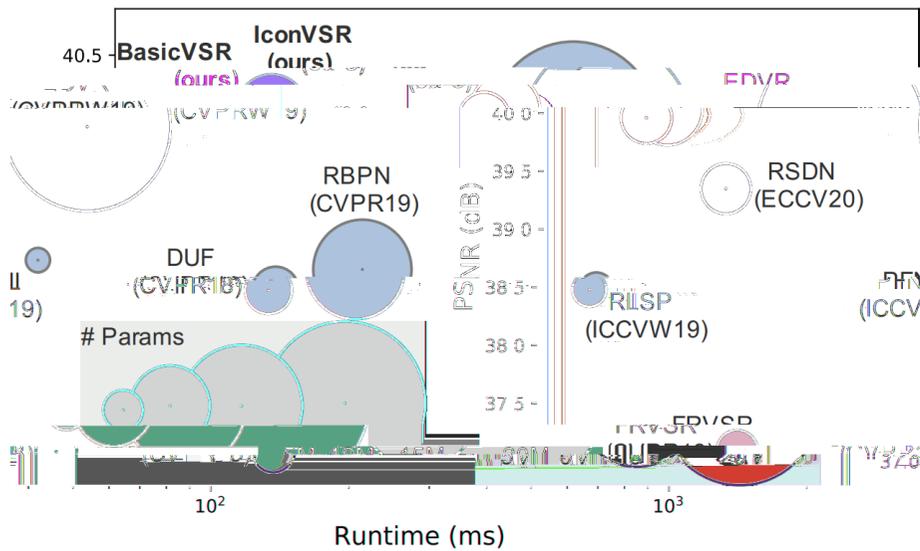
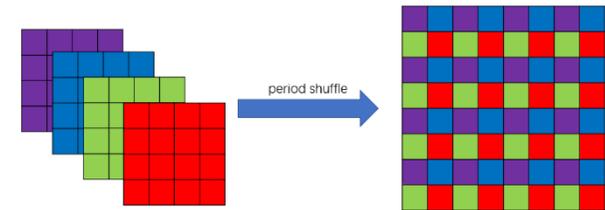
e.g. BCHW=7x 48x 480x 720



e.g. BCHW=7x 3x 1920x 1080



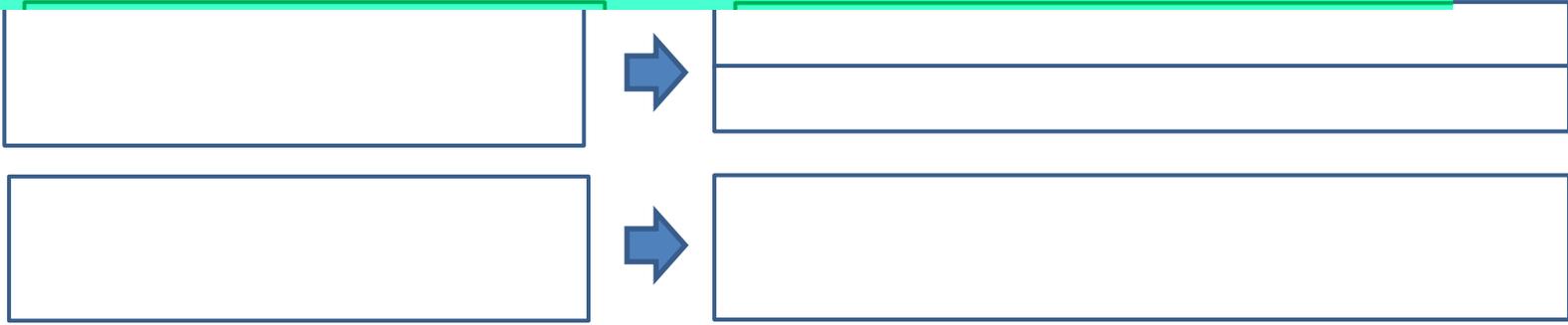
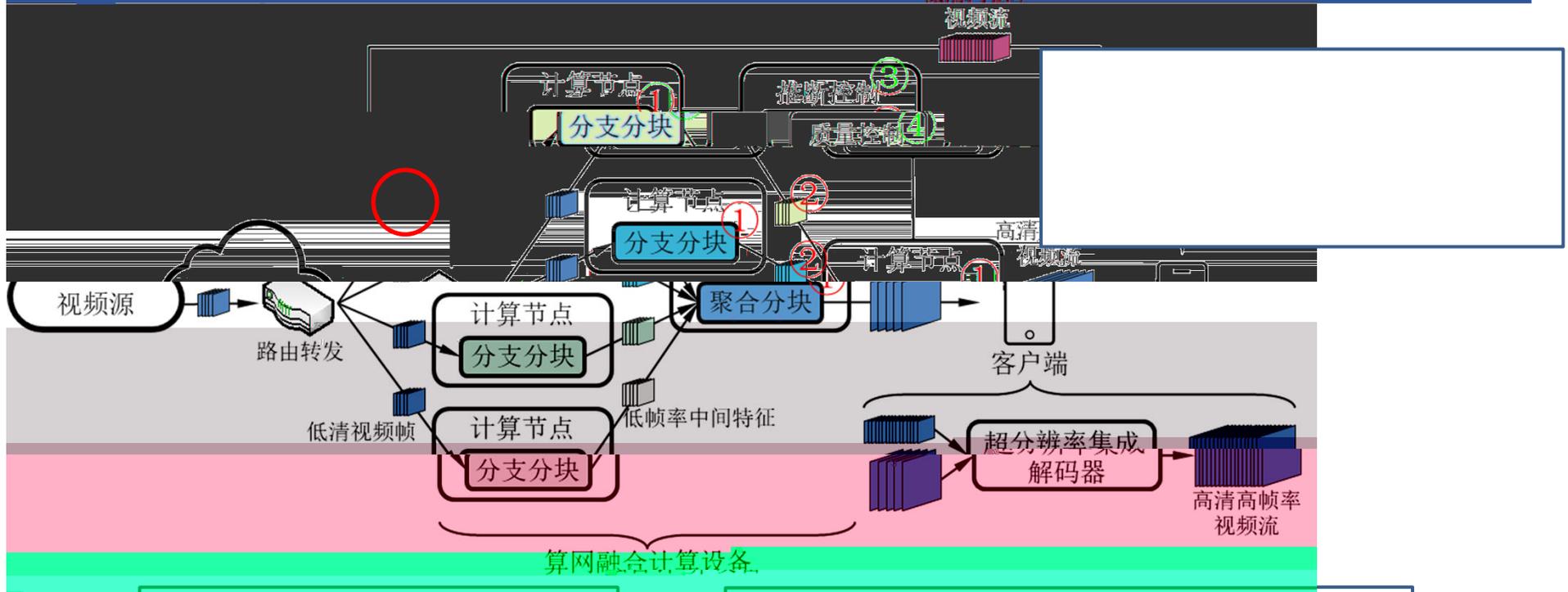
(S SR)
(V SR)





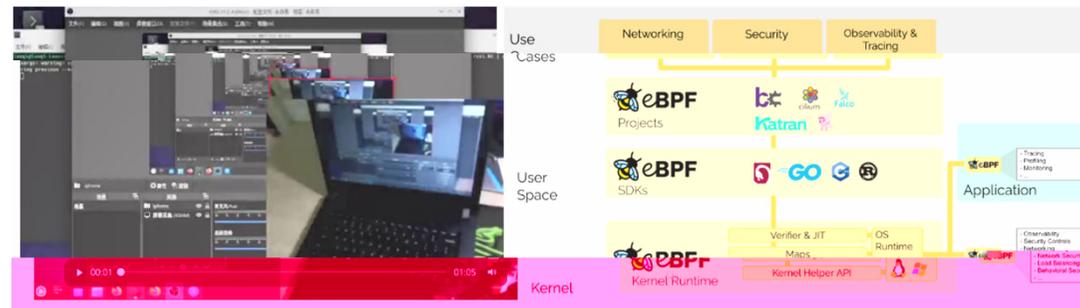
中心

多媒体网络传输技术及应用



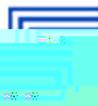


Linux eBPF



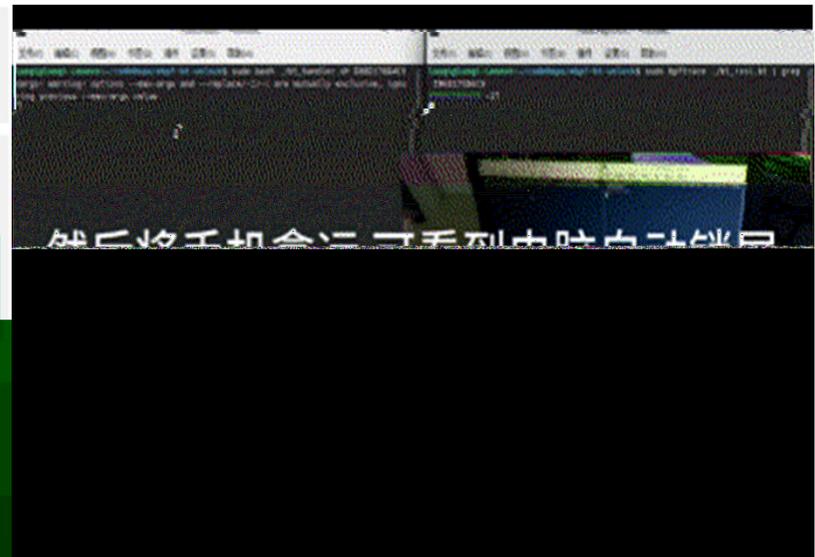
eBPF

Distributed and Optimal RDMA Resource Scheduling in Shared Data Center Networks, Infocom 2020. **CCF A**
 Facilitating Application-aware Bandwidth Allocation in the Cloud with One-step-ahead Traffic Information. **CCF A**
 Last-mile Matters: Mitigating the Tail Latency of Virtualized Networks with Multipath Data Plane. IEEE CLUSTER 2022. **CCF B**



☐ eBPF:

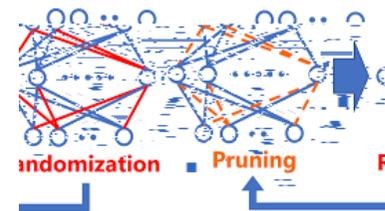
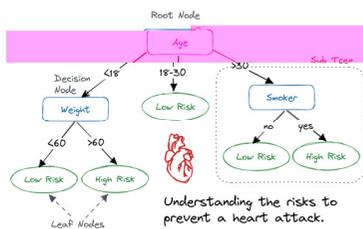
extended Berkeley Packet Filter





eBPF

eBPF

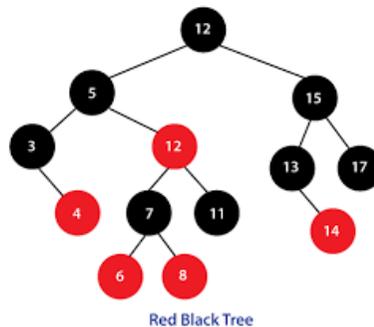


Linux eBPF

■ eBPF

■ memory

dynamic unbounded-loop



```

Algorithm 1: Insertion operation in Heap
Data: B: input array; N: starting index; newValue: new node
Result: Heap tree with the new node
Procedure Insertion(B, N, newValue)
  N = N + 1;
  B[N] = newValue;
  k = N;
  while k > 1 do
    PNode = k/2;
    if B[PNode] < B[k] then
      swap(B[PNode], B[k]);
      k = PNode;
    else
      return;
    end
  end
end

```



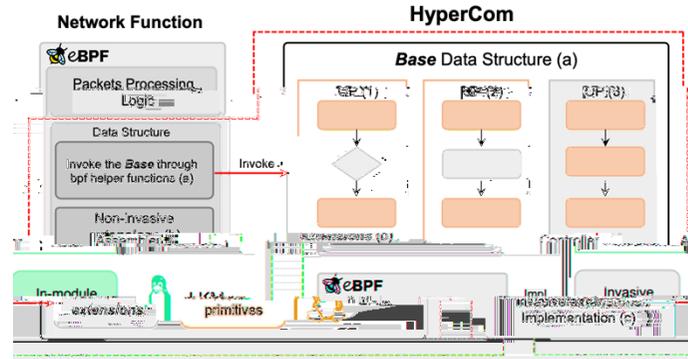
eBPF

HyperCom: Enabling High Performance Complex Data Structure Designs with eBPF for Intelligent Applications

- LKM base part
eBPF extension part (Comprehensiveness)

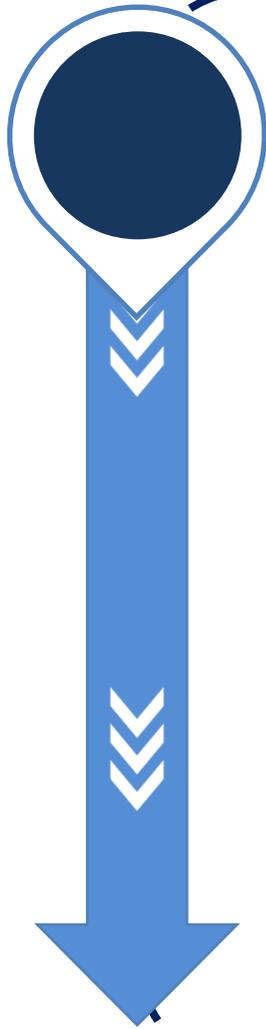
- (eg: SIMD) (Performance)

- (Flexibility)

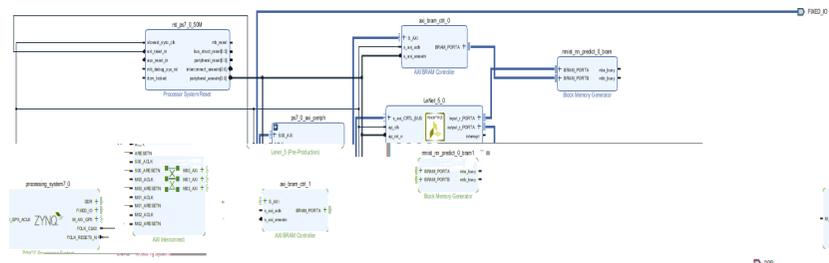


62.3%, 96.8% 89.2%





FPGA



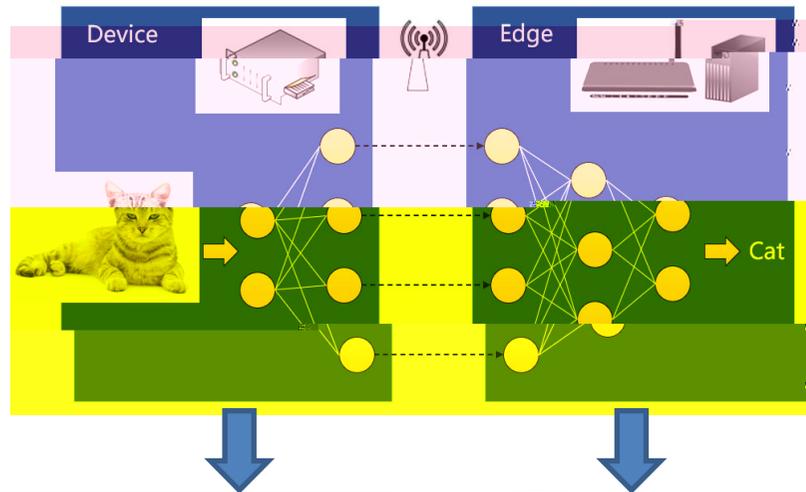
CPU

FPGA





❑ FPGA: Field Programmable Gate Array



FPGA
-
FPGA
-

HLS



The screenshot shows the Vivado IDE interface. The top left displays 'IP时序' (IP Timing) waveforms. The top right shows '基于vivado搭建电路' (Circuit built with Vivado). The bottom left shows code for '基于vitisHLS开发IP核' (IP core developed with Vitis HLS). The bottom right shows '配置ZYNQIP核' (Configure ZYNQIP core) settings.



- Dian Shen, Junzhou Luo, Fang Dong*, Xiaolin Guo, Ciyuan Chen, Kai Wang, John C.S Lui. Enabling Distributed and Optimal RDMA Resource Sharing in Large-scale Data Center Networks: Modeling, Analysis, and Implementation. IEEE/ACM Transactions on Networking, 2023. Accept. **CCF A**
- Bin Yang, Dian Shen, Junxue Zhang, Fang Dong, Junzhou Luo, John C.S. Lui. Towards the Full Extensibility of Multipath TCP with eMPTCP. 2022 IEEE International Conference on Network Protocols (ICNP). **(TH-CPL A)**
- Fang Dong, Huitian Wang, Dian Shen, et al. Multi-exit DNN Inference Acceleration based on Multi-Dimensional Optimization for Edge Intelligence. IEEE Transactions on Mobile Computing, 2022, 10.1109/TMC.2022.3172402. **CCF A**
- Dian Shen, Junzhou Luo, et al. Distributed and Optimal RDMA Resource Scheduling in Shared Data Center Networks, Infocom 2020: 606-615. **(CCF A)**
- Shen Dian, Luo Junzhou, et al. Facilitating Application-aware Bandwidth Allocation in the Cloud with One-step-ahead Traffic Information [J]. IEEE Transactions on Services Computing, 2020, 13(2):381-394. **(CCF A)**

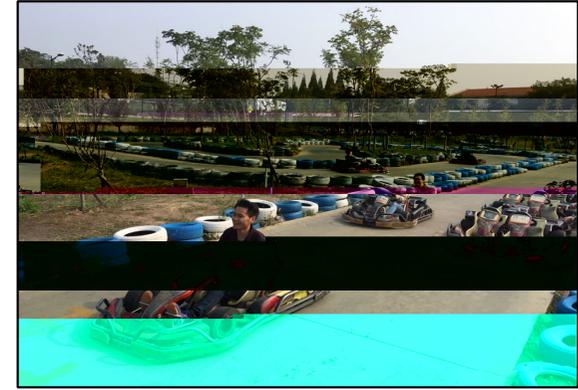
CCF A/B

-
-
-
-

Small and sweet



计算机科学与工程学院
School of Computer Science and Engineering



“

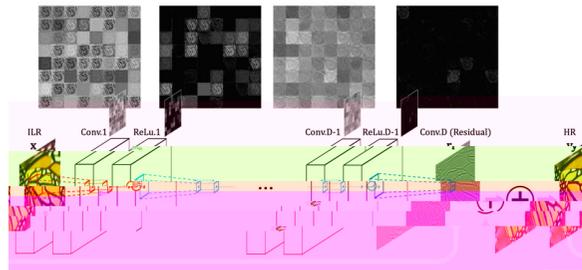
”



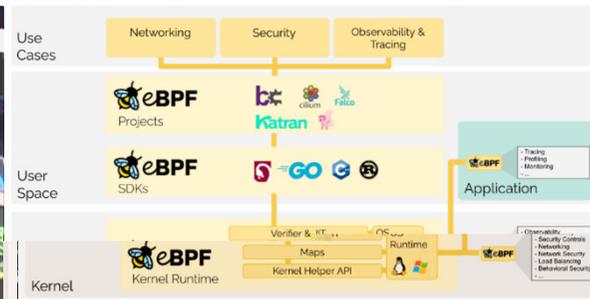
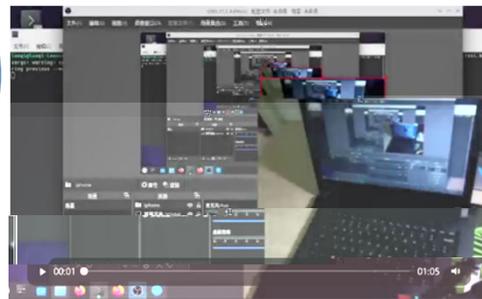
计算机科学与工程学院
School of Computer Science and Engineering



- dshen@seu.edu.cn
- 18100621588



Linux eBPF



FPGA

