Task-oriented communication for edge AI

From "how to communicate" to "what to communicate"

Jun Zhang



MeditCom 2021 tutorial

1

Outline

- Introduction
- Edge Computing and Edge AI
 - Mobile edge computing
 - Edge-assisted inference
- From data-oriented to task-oriented communication
- Task-oriented communication
 - End-to-end design
 - Feature encoding via information bottleneck
 - Communication-computation tradeoff
 - Distributed feature encoding for cooperative inference
- Conclusions

From mobile Internet of mobile intelligence

"When wireless is perfectly applied the whole earth will be converted into a huge brain, which in fact it is, all things being particles of a real and rhythmic whole. We shall be able to communicate with one another instantly, irrespective of distance. Not only this, but through television and telephony we shall see and hear one another as perfectly as though we were face to face, despite intervening distances of thousands of miles; and the instruments through which we shall be able to do this will [fit in a] vest pocket."

The era of Mobile Internet



MeditCom 2021 tutorial -- "Task-oriented communication"





https://marionoioso.com/2018/01/08/from-digital-first-to-ai-first/

Mobile Intelligence







New communication challenges

• Enormous volume of data

- For example, 4TB sensing data/day for autonomous vehicles
- Communication for computing/inference
 - Mobile edge computing, edge AI
- Low-latency communication
 - Millisecond-level latency for safety-critical applications
- Resource-constrained devices
 - Limited onboard computation and communication resources

Edge Computing and Edge AI

From Cloud to Edge

Edge Computing and Edge Al

The Tipping Point



MeditCom 2021 tutorial -- "Task-oriented communication"

Growth of Mobile App Markets



MeditCom 2021 tutorial -- "Task-oriented communication"

Edge Computing and Edge AI

NEED for intensive and low-latency computation









MeditCom 2021 tutorial -- "Task-oriented communication"

Edge Computing and Edge AI

NEED FOR **SPEED**



- VR/AR
 - Latency < 20 ms</p>
 - Avoid cybersickness



- Autonomous Driving
 - For platooning control
 - Latency < 100 ms</p>

MeditCom 2021 tutorial -- "Task-oriented communication"

Old Paradigm for Mobile Computing (I)

- Cloud computing
 - "Internet-based computing that provides shared computer processing resources and data to computers and other devices on demand"



MeditCom 2021 tutorial -- "Task-oriented communication"

Old Paradigm for Mobile Computing (II)

• Mobile cloud computing (MCC)



A new paradigm: Mobile Edge Computing (MEC)



- European Telecommunications Standard Institute (ETSI), 2014
 - MEC "provides IT and cloud-computing capabilities within the Radio Access Network (RAN) in close proximity to mobile subscribers"



MeditCom 2021 tutorial -- "Task-oriented communication"

Edge Computing and Edge AI

Communication Latency



MEC vs. MCC

	Mobile Edge Computing	Mobile Cloud Computing
Hardware	Small-scale data centers	Large-scale data centers
Server location	Co-located with wireless gateways, WiFi routers and BSs	Installed at dedicated buildings
Deployment	Lightweight configuration and planning	Sophisticated configuration and planning
Backhaul Usage	Infrequency use, alleviate congestion	Frequent use, likely to cause congestion
Distance to Users	Tens to hundreds of meters	Across the country boarders

Edge Computing and Edge AI

Application: XR



MeditCom 2021 tutorial -- "Task-oriented communication"

Edge-assisted VR



Fig. 7. FURION system architecture.

Z. Lai, Y. C. Hu, Y. Cui, L. Sun, N. Dai and H. Lee, "Furion: Engineering High-Quality Immersive Virtual Reality on Today's Mobile Devices," in *IEEE Transactions on Mobile Computing*, vol. 19, no. 7, pp. 1586-1602, 1 July 2020

Application: Autonomous driving



Edge-assisted autonomous driving





Edge Computing and Edge AI

Al revolution



Edge Computing and Edge Al

Edge AI is coming



Accelerating AI on the intelligent edge: Microsoft and Qualcomm create vision AI developer kit 000

Posted on May 7, 2018

Edge Computing and Edge Al

From "edge computing" to "edge AI"



References on edge AI

- Z. Zhou, X. Chen, E. Li, L. Zeng, K. Luo and J. Zhang, "Edge Intelligence: Paving the Last Mile of Artificial Intelligence With Edge Computing," *Proceedings of the IEEE*, vol. 107, no. 8, pp. 1738-1762, Aug. 2019.
- X.Wang,Y. Han,V. C. M. Leung, D. Niyato, X.Yan and X. Chen, "Convergence of Edge Computing and Deep Learning: A Comprehensive Survey," *IEEE Communications Surveys & Tutorials*, vol. 22, no. 2, pp. 869-904, Second quarter 2020
- G. Zhu, D. Liu, Y. Du, C. You, **J. Zhang**, and K. Huang, "Toward an intelligent edge: Wireless communication meets machine learning," *IEEE Commun. Mag.*, vol. 58, no. 1, pp. 19–25, Jan. 2020.
- Y. Shi, K. Yang, T. Jiang, **J. Zhang**, and K. B. Letaief, "Communication-efficient edge AI: Algorithms and systems," *IEEE Commun. Surveys Tuts.*, 4th Quart. 2020.
- J. Zhang and K. B. Letaief, "Mobile edge intelligence and computing for the Internet of vehicles," *Proc. IEEE*, vol. 108, no. 2, pp. 246–261, Feb. 2020.
- L.Wang, J. Zhang, J. Chuan, R. Ma, and A. Fei, "Edge intelligence for mission cognitive wireless emergency networks," *IEEE Wireless Commun. Mag.*, vol. 27, no. 4, pp. 103–109, Aug 2020.

Edge-assisted Inference



MeditCom 2021 tutorial -- "Task-oriented communication"

Edge Inference

- Perform low-latency Al inference at the wireless network edge.
- Problems: constrained on-device resources and limited wireless bandwidth.







Aerial Tracking

Edge-assisted inference

Challenge of edge inference – LARGE size of DNN models



https://www.topbots.com/a-brief-history-of-neural-network-architectures/

Edge-assisted inference

Challenge of edge inference – LARGE size of DNN models



https://ai.googleblog.com/2019/05/efficientnet-improving-accuracy-and.html, May 2019

Challenge of edge inference – High energy of DNN models



MeditCom 2021 tutorial -- "Task-oriented communication"

Challenge of Edge Inference – Limited On-device Resources



Challenge of Edge Inference – Limited On-device Resources



Edge-assisted inference

Challenge: How to run powerful deep learning models at resource-constrained devices?
Three different approaches of edge inference



MeditCom 2021 tutorial -- "Task-oriented communication"

Edge-assisted inference

Solution I: Offloading to an edge server

- High communication overhead
 - 256×256 camera with 24 FPS creates around 1.5 MB data per second.
 - 4k camera with 30 FPS creates around 230MB data per second.
- Unstable communication link
 - Wi-Fi: 18 Mbps
 - 4G: 5.85 Mbps
 - 3G: I.I Mbps
- Data privacy



high and a state of the state o



Solution II: On-device inference with model compression

• High compression ratio \rightarrow large performance loss



Han, Song, Huizi Mao, and William J. Dally. "Deep compression: Compressing deep neural networks with pruning, trained quantization and huffman coding." arXiv preprint arXiv:1510.00149 (2015).

Solution III: Device-edge co-inference

• Provide better communication vs. on-device computation tradeoff



MeditCom 2021 tutorial -- "Task-oriented communication"

Edge-assisted inference

Model splitting: An example

Neurosurgeon



Figure 10: Overview of Neurosurgeon. At deployment, Neurosurgeon generates prediction models for each layer type. During runtime, Neurosurgeon predicts each layer's latency/energy cost based on the layer's type and configuration, and selects the best partition point based on various dynamic factors.

Y. Kang, J. Hauswald, C. Gao, A. Rovinski, T. Mudge, J. Mars, and L. Tang. "Neurosurgeon: Collaborative intelligence between the cloud and mobile edge," *ACM SIGPLAN Notices*, 52(4):615–629, 2017.

Edge-assisted inference

Model splitting: Where to split?

- Not all layer are suitable to split
 - Depending on the neural network architecture



MeditCom 2021 tutorial -- "Task-oriented communication"

Model splitting: Data amplification



Fig. 2: Data size at each decoupling point in ResNet, showing that the input data is amplified in in-layer feature maps.

Li H, Hu C, Jiang J, et al. "Jalad: Joint accuracy-and latency-aware deep structure decoupling for edge-cloud execution", 2018 IEEE 24th International Conference on Parallel and Distributed Systems. IEEE, 2018: 671-678.

Model splitting as a computation offloading problem

 Model splitting can be regarded as a computation offloading problem with dependency among sub-tasks.



Fig. 4: A 7-layer DNN model classifies frames of video.



Fig. 5: The inception v4 network represented in layer form.

Fig. 6: Graph representation of inception v4 network.

C. Hu, W. Bao, D. Wang, and F. Liu, "Dynamic Adaptive DNN Surgery for Inference Acceleration on the Edge," in *Proc. IEEE INFOCOM 2019*, 2019, pp. 1423-1431.

Edge-assisted inference

Feature compression

- The intermediate feature is sparse
 - Lossless compression (ZIP, PNG, GIF)
- The neural network (NN) has fault-tolerant property
 - Lossy compression (JPEG, DCT, Quantification)



Feature compression: An example

BottleNet



A. E. Eshratifar, A. Esmaili, and M. Pedram, "Bottlenet: A deep learning architecture for intelligent mobile cloud computing services," arXiv preprint arXiv:1902.01000, 2019.

Feature transmission: Traditional communication

• Separate design of source and channel coding



MeditCom 2021 tutorial -- "Task-oriented communication"

Feature transmission: Joint source-channel coding

- Joint source-channel coding
 - Use learning-based method for effective design
 - The object is to transmit the data reliably \sim



E. Bourtsoulatze, D. B. Kurka, and D. Gndz, "Deep joint source-channel coding for wireless image transmission," in ICASSP 2019, May 2019, pp. 4774-4778.

Edge-assisted inference

The Big Picture



MeditCom 2021 tutorial -- "Task-oriented communication"

Edge-assisted inference

Wishful thinking

- Pseudo formulation
 - **Minimize** Inference Latency (= computation + communication)
 - Subject to On-device computation cost < c Inference error < t

How to make it more tractable?

From "edge computing" to "edge Al"

• Essential idea:

- To provide DNN-based high-performance inference (target) on resourceconstrained devices (constraint) by jointly designing local computation and offloading/communication (methodology)
- Key considerations
 - Communication: informative and concise representation
 - **Computation**: low-complexity feature extractor and encoder

Interplay between communication, computation, and machine learning

A paradigm shift

From data-oriented to task-oriented communication

Digital communications 101



MeditCom 2021 tutorial -- "Task-oriented communication"

Digital communications 101

• Shannon's wisdom



"The fundamental problem of communication is that of reproducing at one point either exactly or approximately a message selected at another point."

The "sematic aspects of communication are irrelevant to the engineering problem. The significant aspect is that the actual message is one selected from a set of possible messages."



C. E. Shannon. A mathematical theory of communication. Bell System Technical Journal, 27:379–423, 623–656, 1948.

Task-oriented communication

55 \cap

A broader view: Three levels of communications

		Shannon's information theory	
	Level A The technical problem	 How accurately can the symbols of communication be transmitted? 	\checkmark
	Level B The semantic problem	• How <i>precisely</i> do the transmitted symbols convey the desired meaning?	
			2
	Level C The effectiveness problem	 How effectively does the received meaning affect conduct in the desired way? 	
ry of ION			,
°	W. Weaver. Recent contributions to a Communication. University of Illinois Pr	the mathematical theory of communication. In C. E. Shannon and W. Weaver, editors, The Mathematical Theory ress, Urbana, 1949.	' of

MeditCom 2021 tutorial -- "Task-oriented communication"

COMMUNI

Three levels of communications



MeditCom 2021 tutorial -- "Task-oriented communication"

Task-oriented communication

A Proposal for the

DARTMOUTH SUMMER RESEARCH PROJECT ON ARTIFICIAL INTELLIGENCE

June 17 - ling. 16

We propose that a 2 month, 10 man study of artificial intelligence be

carried out during the summer of 1956 at Dartmouth College in Hanover, New Hampshire. The study is to proceed on the basis of the conjecture that every aspect of learning or any other feature of intelligence can in principle be so precisely described that a machine can be made to simulate it. An attempt will be made to find how to make machines use language, form abstractions and concepts, solve kinds of problems now reserved for humans, and improve themselves. We think that a significant advance can be made in one or more of these problems if a carefully selected group of scientists work on it together for a summer.

The following are some aspects of the artificial intelligence problem:

1) Automatic Computers

If a machine can do a job, then an automatic calculator can be programmed to simulate the machine. The speeds and memory capacities of present computers may be insufficient to simulate many of the higher functions of the human brain, but the major obstacle is not lack of machine capacity, but our inability to write programs taking full advantage of what we have.

How Can a Computer be Programmed to Use a Language
 It may be speculated that a large part of human thought consists of manipulating words according to rules of reasoning



1956 Dartmouth Conference

Data-oriented vs. task-oriented communication



Challenges and potential solutions of task-oriented communication

• Challenges



Task-oriented communication



Related endeavors in rethinking communication problem

Deep Learning Enabled Semantic Communication Systems

Huiqiang Xie, Zhijin Qin, Member, IEEE, Geoffrey Ye Li, Fellow, IEEE, and Biing-Hwang Juang Life Fellow, IEEE

Semantic Communications in Networked Systems

Elif Uysal, Onur Kaya, Anthony Ephremides, James Gross, Marian Codreanu, Petar Popovski, Mohamad Assaad, Gianluigi Liva, Andrea Munari, Touraj Soleymani, Beatriz Soret and Karl Henrik Johansson*

6G Networks: Beyond Shannon Towards Semantic and Goal-Oriented Communications

Emilio Calvanese Strinati^{a,1,*}, Sergio Barbarossa^b

Task-oriented Communication System Design in

Cyber-Physical Systems:

A Survey on Theory and Applications

Arsham Mostaani, *Student Member, IEEE*, Thang X. Vu, *Member, IEEE*, and Symeon Chatzinotas, *Senior Member, IEEE*

MeditCom 2021 tutorial -- "Task-oriented communication"

62

End-to-end design: BottleNet++, Branchy-GNN

J. Shao, **J. Zhang**, "BottleNet++: An end-to-end approach for feature compression in device-edge coinference systems," *IEEE Int. Conf. Commun. (ICC)* Workshop 2020.

J. Shao, H. Zhang, Y. Mao, and **J. Zhang**, "Branchy-GNN: a device-edge co-inference framework for efficient point cloud processing," *ICASSP 2021*.

Existing studies

• Separately design model splitting, feature compression, and transmission.



MeditCom 2021 tutorial -- "Task-oriented communication"

End-to-end design via deep learning



MeditCom 2021 tutorial -- "Task-oriented communication"

End-to-end design via deep learning



Case study I: BottleNet++ for image classification

BottleNet++

- End-to-End architecture.
- The encoder and decoder are a pair of complementary CNNs.
- The wireless channel is modeled as a non-trainable layer (a transfer function).
- Key insight: Accurate data recovery after transmission is not needed.



Proposed Architecture

Encoder

- Plays the role of feature compression and transmission (joint source-channel coding).
- Key components:
 - A convolutional layer: to reduce the dimension of intermediate feature tensor.
 - A batch normalization layer and a Sigmoid function: to introduce non-linearity, which improve the NN compression capability.
- The channel condition will be an additional input, e.g., variance of the noise or bit erasure rate, for encoder.



Proposed Architecture

Channel

- The wireless channel is modeled as a non-trainable layer.
- AWGN channel: f(x) = x + n $n \sim \mathcal{N}(0, \sigma^2)$
- Binary Erasure Channel (BEC):

$$x \stackrel{Quantify}{\longrightarrow} \hat{x} \stackrel{BEC}{\longrightarrow} \tilde{x} \stackrel{Recovery}{\longrightarrow} f(x)$$



MeditCom 2021 tutorial -- "Task-oriented communication"

Proposed Architecture

Decoder

- Plays the role of joint source-channel decoding.
- Key components:
 - A convolutional layer: to restore the corrupted feature.
 - A batch normalization layer and a Sigmoid function: to introduce non-linearity.



MeditCom 2021 tutorial -- "Task-oriented communication"

Proposed Architecture

Three-step Training

• First Step: train the DNN to reach the desired accuracy, find the model splitting point.



Proposed Architecture

Three-step Training

• Second Step: train the encoder and decoder, fixing other layers.


Proposed Architecture

Three-step Training

• Third Step: fine-tune all layers.



Experiment

- Experiment setup
 - DNN models:VGG16 and ResNet50.
 - Wireless Channel: the AWGN channel and Binary Erasure Channel (BEC).
 - Dataset: CIFAR-100.
 - Task: Image classification.



MeditCom 2021 tutorial -- "Task-oriented communication"



Experiment

- Compression capability comparison
 - Other intermediate feature compression methods.
 - Quantize the intermediate feature and the use the Huffman code to compress it. (Quantification + Huffman)
 - 2. Use JPEG algorithm to compress the intermediate feature. (JPEG)
 - 3. Compress the feature using learning based method and JPEG algorithm (BottleNet)
 - Evaluation criterion
 - Communication overhead (transmitted data size/transmitted latency)
 - On-device computation (approximated by float-point multiplication in convolutional layer)
 - The accuracy degradation < 2%.
 - Definition of peak signal-to-noise ratio (PSNR) in AWGN channel.

$$PSNR = 10 \log_{10} \frac{1}{\sigma^2} \quad (dB)$$

Experiment Result



MeditCom 2021 tutorial -- "Task-oriented communication"

Experiment

• Split the network when the communication overhead less than transmitting PNG image.



On-device Computation	BottleNet++ (ours)	BottleNet	JPEG	Quan.+Huffman
ResNet50 (BEC)	$9.7 imes10^{8}$	1.5×10^{9}	1.6×10^{9}	
ResNet50 (AWGN)	$4.8 imes10^8$	1.4×10^{9}	1.6×10^{9}	
VGG16 (BEC)	$1.4 imes10^9$	3.0×10^9	3.0×10^9	3.0×10^9
VGG16 (AWGN)	$1.4 imes10^9$	3.0×10^{9}	3.0×10^{9}	3.0×10^{9}

2x~3x

MeditCom 2021 tutorial -- "Task-oriented communication"

Experiment Result

Generalization Ability and Robustness Analysis

- **Case I**: The encoder knows the channel condition in both the training and testing processes.
- Case 2: The encoder only knows the channel condition in the training process. In the testing process, the encoder assumes the channel condition to be 15dB in the AWGN channel or 0.125 in the BEC.
- Case 3: The encoder does not know the channel condition in neither the training nor the testing process.



MeditCom 2021 tutorial -- "Task-oriented communication"

78

Case study 2: Branchy-GNN for point cloud processing

- Graph neural networks (GNNs) for point cloud processing.
- Generate edge features for the point cloud to build a graph representation.



Case study 2: Branchy-GNN for point cloud processing

- New challenges
 - More serious data amplification (high-dimensional node feature)
 - Hard-to-compress graph data (complex dependence between adjacent matrix and node/edge features)



S., Martin, and N. Komodakis. "GraphVAE: Towards generation of small graphs using variational autoencoders." In International Conference on Artificial Neural Networks, pp. 412-422. Springer, Cham, 2018.

Solution: Branch structure

- Early exiting from the main branch
- Intermediate feature compression
 - Readout layer:
 - maps the graph representation to a fixed-size vector
 - Joint Source-Channel Coding (JSCC)
 - Directly map the input feature to the channel symbols.
 - Train the coding scheme with channel noise in an end-to-end manner.
- Server-based network



Branchy-GNN framework

- Achieve a better communication-computation tradeoff
 - Flexibly select the exit point according to the edge environments.



MeditCom 2021 tutorial -- "Task-oriented communication"

Experiment

- Setup
 - Point cloud classification based on ModelNet40 dataset
 - 4 exit points in Branchy-GNN
- Baselines:
 - Model splitting
 - Server-only inference
 - Edge-only inference

Experiment

Edge inference speedup



MeditCom 2021 tutorial -- "Task-oriented communication"

Experiment

Robustness of the proposed method

• Train the network at the SNR of 20dB and test it with SNR ranging from 18dB to 25dB.



MeditCom 2021 tutorial -- "Task-oriented communication"

Task-oriented Communication via Information Bottleneck

J. Shao, Y. Mao, and **J. Zhang**, "Learning task-oriented communication for edge inference: An information bottleneck approach," JSAC, to appear. (<u>https://arxiv.org/abs/2102.04170</u>)

Feature encoding via information bottleneck



N. Tishby, F.C. Pereira, and W. Biale. The information bottleneck method. In The 37th annual Allerton Conf. on Communication, Control, and Computing, 1999.

Information bottleneck

Information Bottleneck



preserving "relevant" information vs. finding "compact" representation

- Applications of information bottleneck
 - Understand deep learning
 - Improve generalization performance and robustness to adversarial attack



Variational Feature Encoding (VFE)



feature

feature



Variational Feature Encoding (VFE)

(Alemi et. al., 2016)

 $p_{\phi}(\hat{z}|x)$ is defined by the neural network (encoder)

- $p(\hat{oldsymbol{z}})$ is an unknown marginal distribution
- $q(\hat{m{z}})$ is a log-uniform distribution to approximate the unknown marginal distribution

Some details: Approximated closed-form solution

 $p_{\phi}(\hat{z}|x)$ is a factorized Gaussian distribution

 $q(\hat{oldsymbol{z}})$ is a log-uniform distribution

$$D_{KL}(p_{\phi}(\hat{z}|x)\|q(x)) = \sum_{i=1}^n D_{KL}(p_{\phi}(\hat{z}_i|x)\|q(\hat{z}_i))$$

$$egin{aligned} -D_{KL}(p_{\phi}(\hat{z}_i \mid oldsymbol{x}) \| q(\hat{z}_i)) &= rac{1}{2} \log lpha_i - \mathbb{E}_{\epsilon \sim \mathcal{N}(1, lpha_i)} \log |\epsilon| + \mathrm{C} \ &pprox k_1 S(k_2 + k_3 \log lpha_i) - 0.5 \logig(1 + lpha_i^{-1}ig) + \mathrm{C} \end{aligned}$$

(Molchanov et. al., 2017)

Empirical results



MeditCom 2021 tutorial -- "Task-oriented communication"

Information bottleneck

Experiment

- **Baselines** (data-oriented communication):
 - DeepJSCC (Joint Source-Channel Coding)
 - Learning-based quantization (w/ ideal channel coding)





- With simple encoder/decoder
- Key enablers
 - End-to-end design
 - Information bottleneck
 - Joint source-channel coding

MeditCom 2021 tutorial -- "Task-oriented communication"

Experiment

• VFE method can better distinguish the data from different classes compared with DeepJSCC.



2-dimensional t-SNE embedding of the received feature in the MNIST classification task with PSNR = 20 dB.

MeditCom 2021 tutorial -- "Task-oriented communication"

Variable-length Variational Feature Encoding (VL-VFE)

- To adapt to channel states: variable-length coding
- To reduce signaling overhead, the coding scheme should be consecutive and monotonic



(b) Consecutive activation MeditCom 2021 tutorial -- "Task-oriented communication"

Variable-length Variational Feature Encoding (VL-VFE)



MeditCom 2021 tutorial -- "Task-oriented communication"

Information bottleneck

Variable-length Variational Feature Encoding (VL-VFE)



MeditCom 2021 tutorial -- "Task-oriented communication"

Improve communication-computation tradeoff

J. Shao, **J. Zhang**, "Communication-computation trade-off in resource-constrained edge inference," *IEEE Commun. Mag.*, Dec 2020.

X. Zhang, J. Shao, Y. Mao, and **J. Zhang**, "Communication-computation efficient device-edge coinference via AutoML," *IEEE GLOBECOM 2021*.

Communication-computation tradeoff



MeditCom 2021 tutorial -- "Task-oriented communication"

A three-step framework

- I. Split the network
- 2. Compress the on-device network
- 3. Encode the intermediate feature



MeditCom 2021 tutorial -- "Task-oriented communication"

Step I: Split the network

- Avoid letting the split points transmit multiple data flows.
- Balance the on-device computation and communication overhead.



Step 2: Compress the on-device network

- Incremental network pruning
 - Pruning the weights based on the I-2 norm.
 - Incrementally increase the sparsity ratio: robust pruning process.
 - Structural network pruning: easy for hardware acceleration.



MeditCom 2021 tutorial -- "Task-oriented communication"

Communication-computation tradeoff

Step 3: Encode the intermediate feature

- Learning-based coding scheme:
 - Source coding: map the input to the learned codewords.
 - Joint Source-channel coding: map the input to the channel symbols.



MeditCom 2021 tutorial -- "Task-oriented communication"

Communication-computation tradeoff

Experiment

- Setup
 - ResNet18 for CIFAR10 classification
- Baselines:
 - BottleNet++ (feature compression)
 - 2-step pruning (model compression)
 - Original network

Experiment

• Communication-Computation Tradeoff



MeditCom 2021 tutorial -- "Task-oriented communication"

Experiment

- Real-word test
 - Raspberry Pi 3 as an edge device
 - PC with RTX 2080 Ti as a server



MeditCom 2021 tutorial -- "Task-oriented communication"

Neural architecture optimization: AutoML

- An automated machine learning (AutoML) framework
 - It determines the sparsity level for each on-device DNN layer and the compression ratio for the intermediate feature vector.
 - A deep deterministic policy gradient (DDPG) algorithm is adopted to choose actions.



MeditCom 2021 tutorial -- "Task-oriented communication"

Communication-computation tradeoff

Communication-computation tradeoff



MeditCom 2021 tutorial -- "Task-oriented communication"
Task-oriented communication for cooperative inference

J. Shao, Y. Mao, and **J. Zhang**, "Task-oriented communication for multi-device cooperative edge inference," submitted to IEEE Transactions on Wireless Communications. <u>https://arxiv.org/abs/2109.00172</u>

New Applications: Cooperative Inference

- Multi-device system.
 - Cooperation among multiple devices with distinct views improves sensing capability.



MeditCom 2021 tutorial -- "Task-oriented communication"

Multi-Device Cooperative Edge Inference

- Device-Edge Co-Inference in multi-device systems.
 - Feature extraction (e.g., the VFE).
 - Distributed feature encoding.



MeditCom 2021 tutorial -- "Task-oriented communication"

Multi-Device Cooperative Edge Inference

Design an efficient method that can fully exploit the correlation among multiple features in distributed feature encoding.



Aguerri, Inaki Estella, and Abdellatif Zaidi. "Distributed variational representation learning." IEEE Trans. Pattern Anal. Machine Intell. 120-138, 2019

Multi-Device Cooperative Edge Inference

- Probabilistic modeling with K devices
- Loss functions



MeditCom 2021 tutorial -- "Task-oriented communication"

Distributed Deterministic Information Bottleneck (DDIB)

✤ DIB objective

$$\mathcal{L}_{\text{DIB}}(\beta) := H(Y \mid U_{1:K}) + \beta \sum_{k=1}^{K} [H(Y \mid U_k) + \underbrace{I(Z_k; U_k)}_{\text{Rate}}]$$

$$\stackrel{\text{\bullet DDIB objective}}{\mathcal{L}_{\text{DDIB}}(\beta)} := H(Y \mid U_{1:K}) + \beta \sum_{k=1}^{K} [H(Y \mid U_k) + \frac{R_{\text{bit}}(U_k)}{R_{\text{bit}}(U_k)}]$$

Proposed method: Variational DDIB (VDDIB)

• Using variational inference to estimate the intractable (entropy) terms.

$$egin{aligned} \mathcal{L}_{ ext{DDIB}}(eta) &:= H(Y \mid U_{1:K}) + eta \sum_{k=1}^{K} [H(Y \mid U_k) + R_{ ext{bit}}(U_k)] \ & igwedge \ & igwed \ & igwedge \ & igwed \ & igwedge \ & igwedge$$

MeditCom 2021 tutorial -- "Task-oriented communication"

Minimizing the VDDIB objective may not result in the optimal raterelevance tradeoff due to the approximations

Introduce a selective retransmission (SR) mechanism to further reduce the communication overhead caused by the redundancy among the extracted features.

Selective Retransmission Mechanism

- The edge server selectively activates the edge devices to retransmit their encoded features based on the informativeness of the received features.
- The mechanism consists of a stopping policy and an attention module.

Selective Retransmission Mechanism

Stopping policy

- Each edge device is allowed to transmit the encoded feature with a maximum number of *T* attempts.
- Once the received features are sufficient to output a confident result, the remaining retransmission attempts can be saved.



MeditCom 2021 tutorial -- "Task-oriented communication"

Selective Retransmission Mechanism

- Attention Module
 - Select the most informative features to retransmit based on the attention scores.



VDDIB with Selective Retransmission Mechanism (VDDIB-SR)

• VDDIB-SR loss function

$$\mathcal{L}_{\text{VDDIB}}(\beta; \boldsymbol{\phi}, \boldsymbol{\psi}) := \mathbf{E}_{p_{\theta}(\boldsymbol{z}_{1:K}, \boldsymbol{y})} \{-\log p_{\psi_{0}}(\boldsymbol{y} \mid \boldsymbol{u}_{1:K}) \\ +\beta \left\{ \sum_{k=1}^{K} -\log p_{\psi_{k}}(\boldsymbol{y} \mid \boldsymbol{u}_{k}) + \sum_{k=1}^{K} R_{\text{bit}}(\boldsymbol{u}_{k}) \right\} \}$$

$$\begin{array}{c} \text{Account for } T \\ \text{predictors} \end{array}$$

$$\mathcal{L}_{\text{VDDIB-SR}}(\beta, T; \tilde{\boldsymbol{\phi}}, \tilde{\boldsymbol{\psi}}, \{\boldsymbol{\psi}_{k}\}_{k=1}^{K}) := \mathbf{E}_{p_{\theta}(\boldsymbol{z}_{1:K}, \boldsymbol{y})} \left\{ \frac{1}{T} \sum_{\tau=1}^{T} -\log p_{\tilde{\psi}_{\tau}}(\boldsymbol{y} \mid \{\tilde{\boldsymbol{u}}_{1:K, t}\}_{t=1}^{\tau}) \right\}$$

$$\left. +\beta \left\{ \sum_{k=1}^{K} -\log p_{\psi_{k}}(\boldsymbol{y} \mid \boldsymbol{u}_{k}) + \sum_{k=1}^{K} \sum_{t=1}^{T} R_{\text{bit}}(\tilde{\boldsymbol{u}}_{k, t}) \right\} \right\}$$

MeditCom 2021 tutorial -- "Task-oriented communication"

Performance Evaluation

• Cooperative inference tasks



Two-view MNIST classification



Twelve-view Shape Recognition on ModelNet40 dataset

			VINIST classification		
			2	$R_{ m sum}$	
			6 bits	10 bits	14 bits
**	The accuracy of the cooperative	NN-REG	95.93%	97.49%	97.78%
	tasks under different hit	NN-GBI	96.62%	97.79%	98.02%
		eSAFS	96.97%	97.87%	98.05%
	constraints.	CAFS	94.14%	97.43%	97.42%
•	Data-oriented communication	VDDIB (ours)	97.08%	97.82%	98.06%
Ť		VDDIB-SR (T=2) (ours)	97.13%	98.13%	98.22%
	leads to	Г			
	 I.3 kbits overhead with 98.6% accuracy in the MNIST 	Shape Recognition			
		$R_{ m sum}$			
	classification task.		120 bits	240 bits	360 bits
	 I 20 KB overhead with 92% 	NN-REG	87.50%	88.25%	89.03%
	accuracy in the shape recognition	NN-GBI*	88.82%	_	_
		eSAFS	85.88%	87.87%	89.50%
	task.	CAFS	86.75%	89.56%	90.67%
		VDDIB (ours)	89.25%	90.03%	90.75%
		VDDIB-SR (T=2) (ours)	90.25%	91.31%	91.62%

Performance Evaluation

* The GBI quantization algorithm is computationally prohibitive when the number of bits is too large.

MeditCom 2021 tutorial -- "Task-oriented communication"

122

Ablation Study

- ✤ Impact of the maximum transmission attempts T.
 - the performance of the VDDIB-SR method improves with T.



MeditCom 2021 tutorial -- "Task-oriented communication"

Ablation Study

Impact of the attention module

 We propose a baseline method that removes the attention module denoted as VDDIB-Cascade for comparison.



MeditCom 2021 tutorial -- "Task-oriented communication"

Ablation Study

- Impact of the attention module
 - We propose a baseline method that removes the attention module denoted as VDDIB-Cascade for comparison.
 - The attention module can schedule the most important features to retransmit.

	accuracy	Bits
VDDIB-SR	91.25%	216
VDDIB-Cascade	90.88%	240





MeditCom 2021 tutorial -- "Task-oriented communication"



Conclusions

- A paradigm shift in communication system design
 - From data-oriented communication to task-oriented communication
 - From "how to communicate" to "what to communicate"
- Enabling techniques
 - End-to-end design via deep learning
 - Information bottleneck principle (with variational approximation)
 - Neural architecture design
- New tradeoffs
 - Communication-accuracy tradeoff
 - Communication-computation tradeoff

An invitation to this exciting topic

• Related research areas

- Deep learning
- Wireless communications
- Information theory
- Distributed computing systems

Application scenarios

- Video analytics
- AR/VR
- Autonomous driving
- Smart drones
- ...

References

- J. Shao, J. Zhang, "BottleNet++: An end-to-end approach for feature compression in device-edge co-inference systems," IEEE Int. Conf. Commun. (ICC) Workshop on Edge Machine Learning for 5G Mobile Networks and Beyond, Jun. 2020.
- J. Shao, **J. Zhang**, "Communication-computation trade-off in resource-constrained edge inference," *IEEE Commun. Mag.*, Dec 2020.
- J. Shao, H. Zhang, Y. Mao, and **J. Zhang**, "Branchy-GNN: a device-edge co-inference framework for efficient point cloud processing," *ICASSP 2021*.
- X. Zhang, J. Shao, Y. Mao, and **J. Zhang**, "Communication-computation efficient device-edge coinference via AutoML," *IEEE GLOBECOM 2021*.
- J. Shao, Y. Mao, J. Zhang, "Learning task-oriented communication for edge inference: An information bottleneck approach," JSAC, to appear. (<u>https://arxiv.org/abs/2102.04170</u>)
- J. Shao, Y. Mao, and **J. Zhang**, "Task-oriented communication for multi-device cooperative edge inference," submitted to *IEEE Trans*. *Wireless Communications*. (<u>https://arxiv.org/abs/2109.00172</u>)

Thank you!

• For more details

http://www.eie.polyu.edu.hk/~jeiezhang/ eejzhang@ust.hk