# Tutorial Agenda

| Morning Session | | |
|---|---|---|
| 9:00 - 9:15 | **Opening Remarks** | Lijuan Wang |
| 9:15 - 10:00 | **Overview of Image-Text Pre-training** | Jianfeng Wang |
| 10:00 - 10:15 | **Coffee Break & QA** | |
| 10:15 - 11:00 | **Unified Image-Text Modeling** | Zhengyuan Yang |
| 11:00 - 11:45 | **Advanced Topics in Image-Text Pre-training** | Zhe Gan |
| 11:45 – 12:00 | **Q & A** | |
| Afternoon Session | | |
| 13:00 - 13:30 | **Overview of Video-Text Pre-training** | Kevin Lin |
| 13:30 - 14:00 | **Learning from Multi-channel Videos: Methods and Benchmarks** | Linjie Li |
| 14:00 - 14:30 | **Advanced Topics in Video-Text Pre-training** | Chung-Ching Lin |
| 14:30 - 14:45 | **Coffee Break & QA** | |
| 14:45 - 15:15 | **VLP for Image Classification** | Jianwei Yang |
| 15:15 - 15:45 | **VLP for Object Detection** | Pengchuan Zhang |
| 15:45 - 16:15 | **Benchmarks for Computer Vision in the Wild** | Chunyuan Li |
| 16:15 - 17:00 | **VLP for Text-to-Image Synthesis** | Chenfei Wu |
| 17:00 - 17:15 | **Q & A** | |

**Zhe Gan**
Microsoft

**Pengchuan Zhang**
Meta

**Zhengyuan Yang**
Microsoft

**Kevin Lin**
Microsoft

**Linjie Li**
Microsoft

**Chunyuan Li**
Microsoft

**Jianfeng Wang**
Microsoft

**Jianwei Yang**
Microsoft

**Chung-Ching Lin**
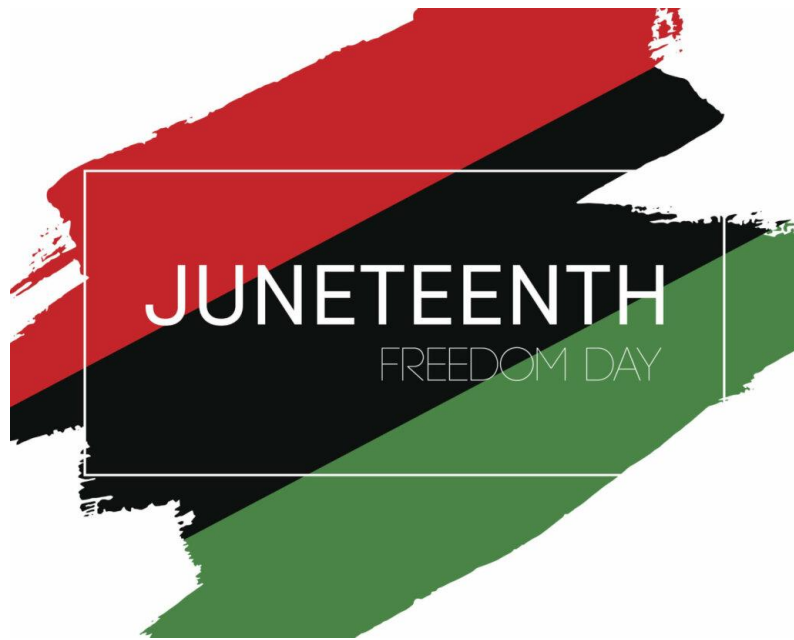Microsoft

**Chenfei Wu**
Microsoft

**Lijuan Wang**
Microsoft

**Zicheng Liu**
Microsoft

**Jianfeng Gao**
Microsoft

Tutorial website: https://vlp-tutorial.github.io/

# Juneteenth Holiday & Father's day



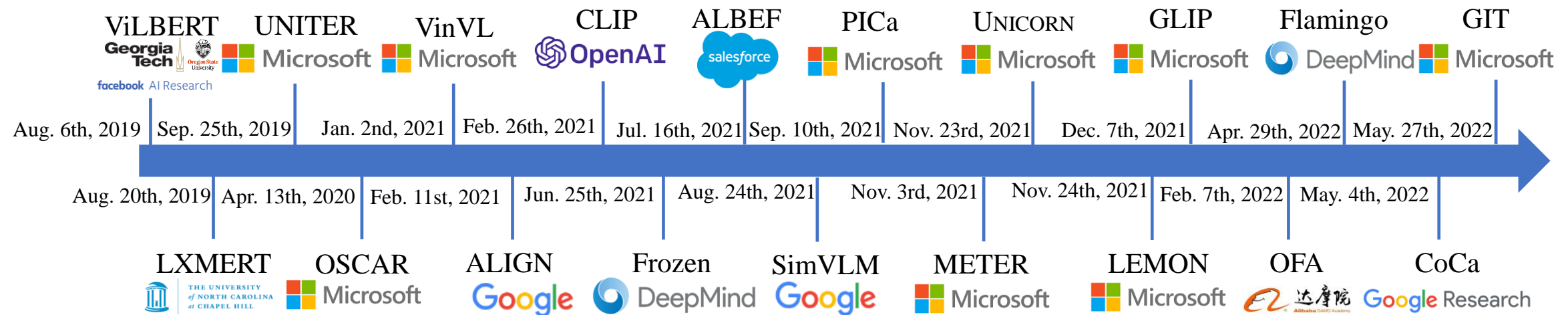https://cvpr2022.thecvf.com/recognizing-juneteenth

# Evolution of Image-Text Pre-training

- All these early models depend on pre-trained object detectors to extract visual features offline.
- Newer end-to-end VLP achieves stronger performance with model and data scaling.
- Upscaled VLP models demonstrate new capabilities such as in-context learning and multimodal few shots.
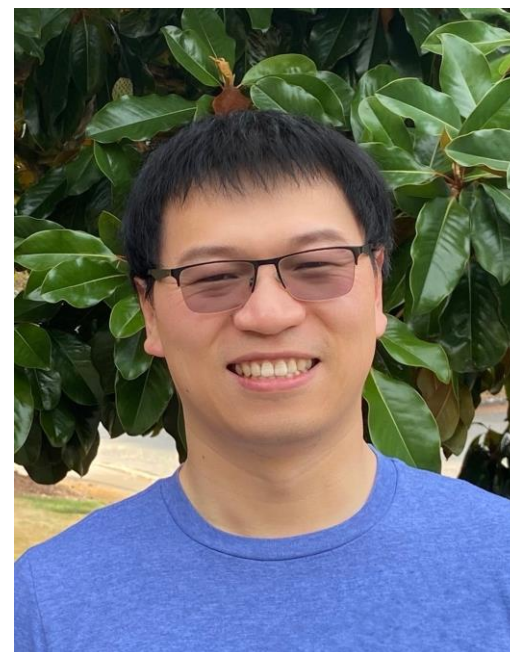
# Session 1: Overview of Image-Text Pre-training

Time:

**10:15 – 11:00 AM**

Presenter:

**Jianfeng Wang (Microsoft)**

Dr. Jianfeng Wang is a Principal Researcher at Microsoft Azure AI. His research interest is computer vision and vision-language intelligence. He received his B. Eng. degree and PhD degree from the University of Science and Technology of China in 2010 and 2015, respectively.

# Session 2: Unified Image-Text Modeling

Time:

**10:15 – 11:00 AM**

Presenter:

**Zhengyuan Yang (Microsoft)**

Dr. Zhengyuan Yang is a Researcher at Microsoft Azure AI. He received the Ph.D. degree in Computer Science at the University of Rochester (UR) in 2021, and the BE degree in electrical engineering from the University of Science and Technology of China (USTC) in 2016. His research interests are vision+language and multi-modal learning.

# Session 3: Advanced Topics in Image-Text Pre-training

Time:

**11:00 – 11:45 AM**

Presenter:

**Zhe Gan (Microsoft)**

Dr. Zhe Gan is a Principal Researcher at Microsoft Azure AI, primarily working on Vision-and-Language Multimodal Intelligence. He received the PhD degree from Duke University in 2018. Before that, he received the Master's and Bachelor's degree from Peking University in 2013 and 2010, respectively. Together with his co-authors, he has received the Best Student Paper Honorable Mention Awards at CVPR 2021 and WACV 2021, respectively.

# Evolution of Video-Text Pre-training

- More advanced models trained end2end on up to 1B video-text pairs with applications to video-text/video-only or even image-text tasks
- Learning from multi-channel videos, leveraging information from vision + language + audio
- Unified architecture across modalities and tasks

# Session 4: Overview of Video-Text Pre-training

Time:

**13:00 – 13:30 PM**

Presenter:

**Kevin Lin (Microsoft)**

Dr. Kevin Lin is an Applied Scientist at Microsoft Azure AI, working on computer vision and vision-language multimodal intelligence. Prior to Microsoft, he obtained his PhD from the University of Washington in 2020, and his MS from National Taiwan University in 2014.

# Session 5: Learning from Multi-channel Videos: Methods and Benchmarks

Time:

**13:30 – 14:00 PM**

Presenter:

**Linjie Li (Microsoft)**

Linjie Li is a Senior Researcher at Microsoft Azure AI. She obtained her Master's degree in computer science from Purdue University in 2018. Her current research interests include Vision-and-Language pre-training and self-supervised learning. Linjie (co-)organizes multiple tutorials on vision+language at CVPR 2020/2021/2022. Her recent work ClipBERT for video+language pre-training is nominated for the Best Student Paper Award at CVPR 2021.

# Session 6: Advanced Topics in Video-Text Pre-training
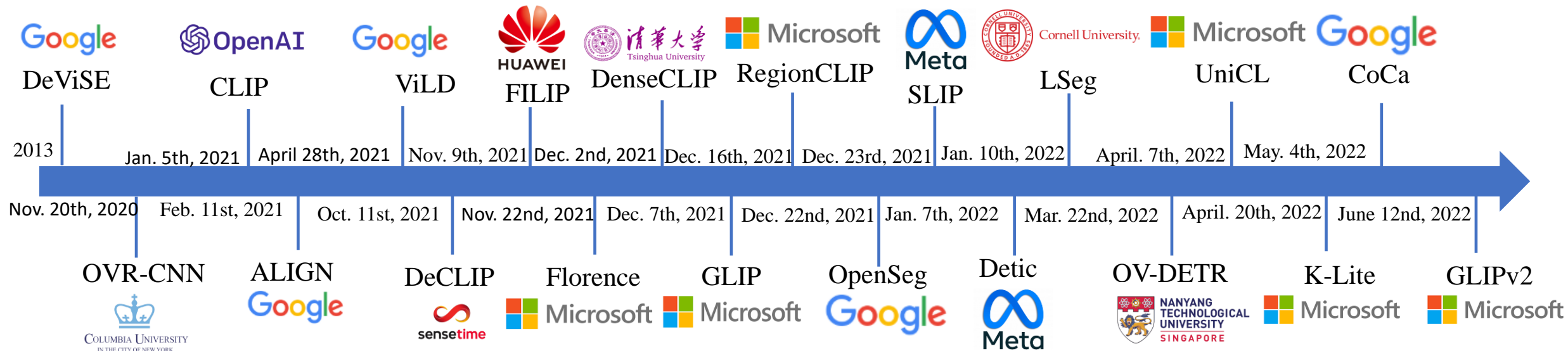
Time:

**13:00 – 13:30 PM**

Presenter:

**Chung-Ching Lin (Microsoft)**

Dr. Chung-Ching Lin is a Principal Researcher at Microsoft Azure AI. Before joining Microsoft, Chung-Ching was a Research Staff Member with IBM T.J. Watson Research, where he began his career after earning his Ph.D. from Georgia Institute of Technology. His research interests include computer vision and machine learning, with specific interests in video understanding, representation learning, and vision and language.

# Evolution of VLP for Vision

- Vision has gradually transit from canonical classification problem to a vision-language alignment problem.
- Contrastive vision-language pretraining supports both representation learning and zero-shot/few-shot transferring learning.
- Vision tasks at image-level, region-level and pixel-level can be unified into a single vision-language interface using contrastive learning and/or generative pretraining.
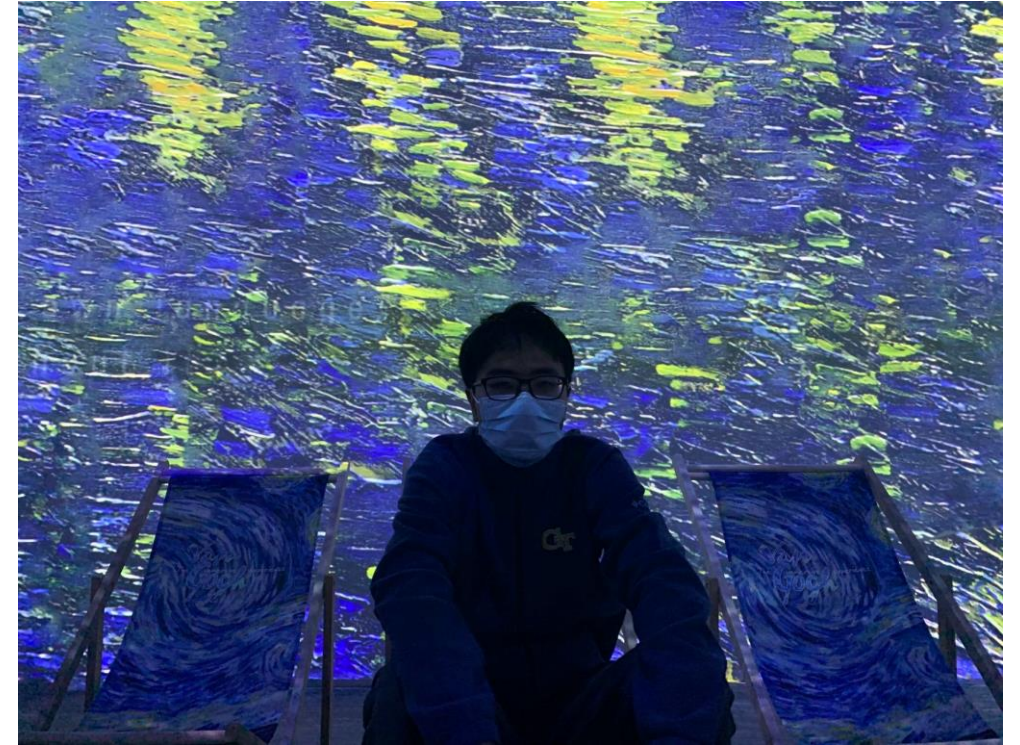
# Session 7: VLP for Image Classification

Time:

**14:45 – 15:15 PM**

Presenter:

**Jianwei Yang (Microsoft)**

Dr. Jianwei Yang is a Senior Researcher at Microsoft Research. His main research interests span in computer vision, vision-language and embodied AI. He obtained his Ph.D. in CS at Georgia Institute of Technology.

# Session 8: VLP for Object Detection

Time:

**15:15 – 15:45 PM**

Presenter:

**Pengchuan Zhang (Meta AI)**

Dr. Pengchuan Zhang is currently an AI research scientist at Meta AI and an affiliate assistant professor at University of Washington. Before joining Meta, he was a Principal Researcher at Microsoft Research. His research interests are mainly in the areas of computer vision, multimodal intelligence, and theoretical foundations for deep learning.

# Session 9: Benchmarks for Computer Vision in the Wild

Time:

**15:45 – 16:15 PM**

Presenter:

**Chunyuan Li (Microsoft)**

Dr. Chunyuan Li is currently a Principal Researcher at Microsoft Research. His recent research focuses on large-scale pre-training in computer vision and natural language processing. He obtained his PhD in machine learning at Duke University.

# Session 10: VLP for Text-to-Image Synthesis

Time:

**16:15 – 17:00 PM**

Presenter:

**Chenfei Wu (Microsoft)**

Dr. Chenfei Wu is currently a Senior Researcher at Microsoft Research. His research interests include visual-language understanding and generation. Recently, he focuses on language guided high-quality image/video generation.

# Tutorial Agenda

| Morning Session | | |
|---|---|---|
| 9:00 - 9:15 | **Opening Remarks** | Lijuan Wang |
| 9:15 - 10:00 | **Overview of Image-Text Pre-training** | Jianfeng Wang |
| 10:00 - 10:15 | **Coffee Break & QA** | |
| 10:15 - 11:00 | **Unified Image-Text Modeling** | Zhengyuan Yang |
| 11:00 - 11:45 | **Advanced Topics in Image-Text Pre-training** | Zhe Gan |
| 11:45 – 12:00 | **Q & A** | |
| Afternoon Session | | |
| 13:00 - 13: 30 | **Overview of Video-Text Pre-training** | Kevin Lin |
| 13:30 - 14:00 | **Learning from Multi-channel Videos: Methods and Benchmarks** | Linjie Li |
| 14:00 - 14: 30 | **Advanced Topics in Video-Text Pre-training** | Chung-Ching Lin |
| 14:30 - 14:45 | **Coffee Break & QA** | |
| 14:45 - 15: 15 | **VLP for Image Classification** | Jianwei Yang |
| 15:15 - 15:45 | **VLP for Object Detection** | Pengchuan Zhang |
| 15:45 - 16:15 | **Benchmarks for Computer Vision in the Wild** | Chunyuan Li |
| 16:15 - 17:00 | **VLP for Text-to-Image Synthesis** | Chenfei Wu |
| 17:00 - 17:15 | **Q & A** | |



**Zhe Gan**
Microsoft

**Pengchuan Zhang**
Meta

**Zhengyuan Yang**
Microsoft

**Kevin Lin**
Microsoft

**Linjie Li**
Microsoft

**Chunyuan Li**
Microsoft

**Jianfeng Wang**
Microsoft

**Jianwei Yang**
Microsoft

**Chung-Ching Lin**
Microsoft

**Chenfei Wu**
Microsoft

**Lijuan Wang**
Microsoft

**Zicheng Liu**
Microsoft

**Jianfeng Gao**
Microsoft

Tutorial website: https://vlp-tutorial.github.io/