# EDGE INTELLIGENCE OPENING

AAAI-EIW-III

https://eiw2024.github.io/

# HISTORY

Home   Scientific Program   Invited Speakers   Challenge   Committees   Submission   Registration   Venue   Book of Abstracts



Edge Intelligence Workshop 2022

19 - 20 September, 2022 - Montréal, Québec, Canada

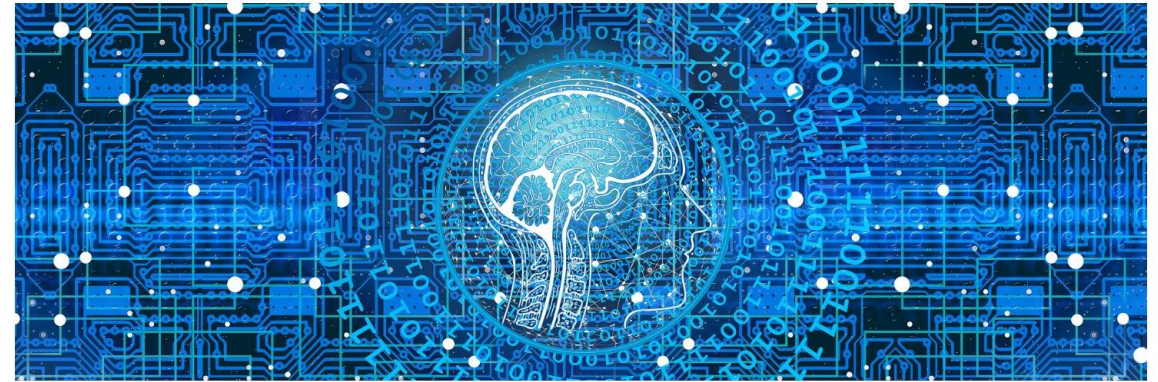## Edge Intelligence Workshop 2022

There is a growing interest towards moving intelligent applications to edge devices, given their advantages such as increased privacy or lower network latency compared to the cloud. However, deploying artificial intelligence systems, which are growing in size, on resource constrained edge devices poses various challenges.
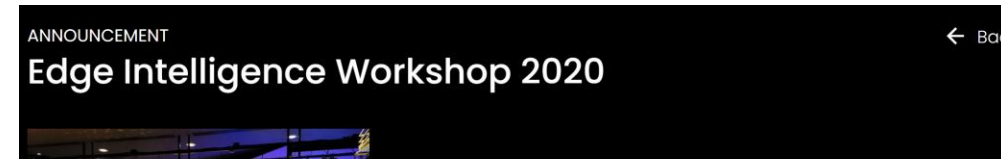
https://eiw2022.github.io/

https://eiw2022.github.io/assets/Proceedings.pdf

https://eiw2024.github.io/

The first series of the Edge Intelligence Workshop was held successfully in 2-3 March 2020 with more than 100 participants from academia and industry. read more here

https://www.gerad.ca/colloques/EdgeIntelligence2020/

https://www.gerad.ca/en/papers/G-2020-23

ANNOUNCEMENT
Edge Intelligence Workshop 2020

← Bac

share

Andrea Lodi & Yoshua Bengio, Edge Intelligence Workshop, March 2-3, 2020

The first series of the Edge Intelligence Workshop was held successfully in 2-3 March 2020 with more than 100 participants from academia and industry.

There has been an emerging interest in edge implementation of deep neural networks, but this direction has no specific scientific event dedicated to it, and this workshop aimed at filling this gap. Edge intelligence is a highly promising area in AI, which is identified as one of the top 10 breakthrough technologies (also known as TinyAI) in 2020 by MIT technology review. This workshop was the first major event in Canada dedicated to this topic.

# WEBSITE AND SLIDES

https://eiw2024.github.io/

EIW 2024    Program    Speakers    Organizers    Scientific Committee    Submission    Venue

## Edge Intelligence Workshop
### Workshop at AAAI 2024
Vancouver - Canada

**The poster session is in room 301 (third floor).**
**Other sessions are held in room 210 (second floor).**

The Edge Intelligence Workshop 2024 will focus on the edge deployment of large language and vision models; and how to make them more efficient in terms of **Data**, **Model**, **Training**, and **Inference** specially on edge devices.

This is an interdisciplinary research topic that covers the theory, hardware, and software aspects of AI models, targeting large language and vision models.

The workshop is part of the **38th Annual AAAI Conference on Artificial Intelligence** and will be held in **Vancouver, Canada**.

**Submission Page**

# ORGANIZERS

**Warren Gross**

McGill University

**Vahid Partovi Nia**

Polytechnique Montreal, Huawei Noah's Ark Lab

**Andrea Lodi**

Cornell Tech

**Shahrokh Valaee**

University of Toronto

**Melika Payvand**

UZH and ETH Zurich

**Mehdi Rezagholizadeh**

Huawei Noah's Ark Lab

**Habib Hajimolahoseini**

Huawei Toronto Research Centre

**Mouloud Belbahri**

Layer6AI TDBank

**Mohammadreza Tayaranian**

McGill University

**Yuanhao Yu**

Huawei Noah's Ark Lab

**Ali Edalati**

Huawei Noah's Ark Lab

# SCIENTIFIC COMMITTEE

| | |
|---|---|
| Abbas Ghaddar | Huawei Noah's Ark Lab |
| Alireza Ghaffari | Huawei Noah's Ark Lab |
| Ali Edalati | Huawei Noah's Ark Lab |
| Didier Chételat | Polytechnique Montreal |
| Ghouthi Boukli-Hacene | Sony |
| Gonçalo Mordido | MILA |
| Erfan Seyedsalehi | Huawei Noah's Ark Lab |
| Ehsan Kamalloo | University of Alberta |
| Hassan Mozafari | McGill University |
| Khalil Bibi | Huawei Noah's Ark Lab |
| Marzieh Tahaei | Huawei Noah's Ark Lab |
| Michael Metel | Huawei Noah's Ark Lab |
| Mohammadreza Tayaranian | McGill University |

| | |
|---|---|
| Mojtaba Valipour | University of Waterloo |
| Peng Lu | Université de Montreal |
| Sharareh Younesian | Huawei Noah's Ark Lab |
| Ramchalam Ramakrishnan | Qualcomm |
| Ritam Haldar | Qualcomm |
| Vanessa Courville | Huawei Noah's Ark Lab |
| Walid Ahmed | Huawei Toronto Research Centre |
| Xinlin Li | Huawei Noah's Ark Lab |
| Zhixiang Chi | Huawei Noah's Ark Lab |

# KEYNOTE SPEAKERS

**Diana Marculescu**

University of Texas at Austin

**Sarath Chandar**

Polytechnique Montreal, MILA

**Lukas Cavigelli**

Huawei Zurich Research Center

**Pascal Poupart**

University of Waterloo

**Di Niu**

University of Alberta

| Start Time | End Time | Program | Speaker |
|---|---|---|---|
| 08:45 | 09:00 | Opening Remarks | |
| 09:00 | 09:30 | Invited Talk | Diana Marculescu |
| 09:30 | 10:00 | Nominated Papers | Diana Marculescu, Mehdi Rezagholizadeh |
| 10:00 | 10:30 | Invited Talk | Di Niu |
| 10:30 | 11:00 | Coffee Break | |
| 11:00 | 11:30 | Invited Talk | Pascal Poupart |
| 11:30 | 12:30 | Poster Discussion & Break | |
| 12:30 | 14:00 | Lunch | |
| 14:00 | 15:00 | Panel Discussion | |
| 15:00 | 15:30 | Invited Talk | Sarath Chandar |
| 15:30 | 16:00 | Coffee Break | |
| 16:00 | 16:30 | Invited Talk | Lukas Cavigelli |
| 16:30 | 16:45 | Closing Remarks | |

SCHEDULE

https://eiw2024.github.io/

# POSTERS 11:30-12:30 THEN LUNCH

- Room 301
- Panel numbers 62-76
- Lunch 12:30-14:00

# PANEL DISCUSSION AFTER LUNCH
# 14:00-15:00

**Sarath Chandar**

Polytechnique Montreal, MILA

**Lukas Cavigelli**

Huawei Zurich Research Center

**Mehdi Rezagholizadeh**

Huawei Noah's Ark Lab

**Habib Hajimolahoseini**

Huawei Toronto Research Centre

**Pascal Poupart**

University of Waterloo

**Di Niu**

University of Alberta

# CLOSING REMARKS
# 16:30-16:45

Prompt:

People watching a presentation about artificial intelligence while being sad that the conference is ending in Vancouver

# EDGE INTELLIGENCE CLOSING

AAAI-EIW-III

11

https://eiw2024.github.io/

# ACCEPTED PAPERS

https://eiw2024.github.io/

EIW 2024     Program   Speakers   Organizers   Scientific Committee   Submission   Venue

Accepted papers

## Edge Intelligence Workshop
### Workshop at AAAI 2024
Vancouver - Canada

The poster session is in room 301 (third floor).
Other sessions are held in room 210 (second floor).

The Edge Intelligence Workshop 2024 will focus on the edge deployment of large language and vision models; and how to make them more efficient in terms of **Data**, **Model**, **Training**, and **Inference** specially on edge devices.

This is an interdisciplinary research topic that covers the theory, hardware, and software aspects of AI models, targeting large language and vision models.

The workshop is part of the **38th Annual AAAI Conference on Artificial Intelligence** and will be held in **Vancouver, Canada**.

Submission Page

# ORGANIZERS

**Warren Gross**

McGill University

**Vahid Partovi Nia**

Polytechnique Montreal, Huawei Noah's Ark Lab

**Andrea Lodi**

Cornell Tech

**Shahrokh Valaee**

University of Toronto

**Melika Payvand**

UZH and ETH Zurich

**Mehdi Rezagholizadeh**

Huawei Noah's Ark Lab

**Habib Hajimolahoseini**

Huawei Toronto Research Centre

**Mouloud Belbahri**

Layer6AI TDBank

**Mohammadreza Tayaranian**

McGill University

**Yuanhao Yu**

Huawei Noah's Ark Lab

**Ali Edalati**

Huawei Noah's Ark Lab

# KEYNOTE SPEAKERS

**Diana Marculescu**
University of Texas at Austin

**Sarath Chandar**
Polytechnique Montreal, MILA

**Lukas Cavigelli**
Huawei Zurich Research Center

**Pascal Poupart**
University of Waterloo

**Di Niu**
University of Alberta

# NOMINATED PAPERS



## SupMAE

- Feng Liang (The University of Texas at Austin)
- Yangguang Li (SenseTime Group Limited)
- Diana Marculescu (The University of Texas at Austin)



## QDyLoRA

- Hossein Rajabzadeh (University of Waterloo)
- Mojtaba Valipour (University of Waterloo)
- Marzieh Tahaei (Huawei Noah's Ark Lab)
- Hyock Ju Kwon (University of Waterloo)
- Ali Ghodsi  (University of Waterloo)
- Boxing Chen   (Huawei Noah's Ark Lab)
- Mehdi Rezagholizadeh (Huawei Noah's Ark Lab)

# NEXT EIW 2026 (BE INVOLVED WITH US)



Prompt:


edge intelligence workshop
happening in an unknown location in
2026

Thank you

SEE YOU IN
TWO YEARS
EIW 2026