# POSTER SESSION: GENERAL IDEA

Similar to a normal presentation

Key differences:
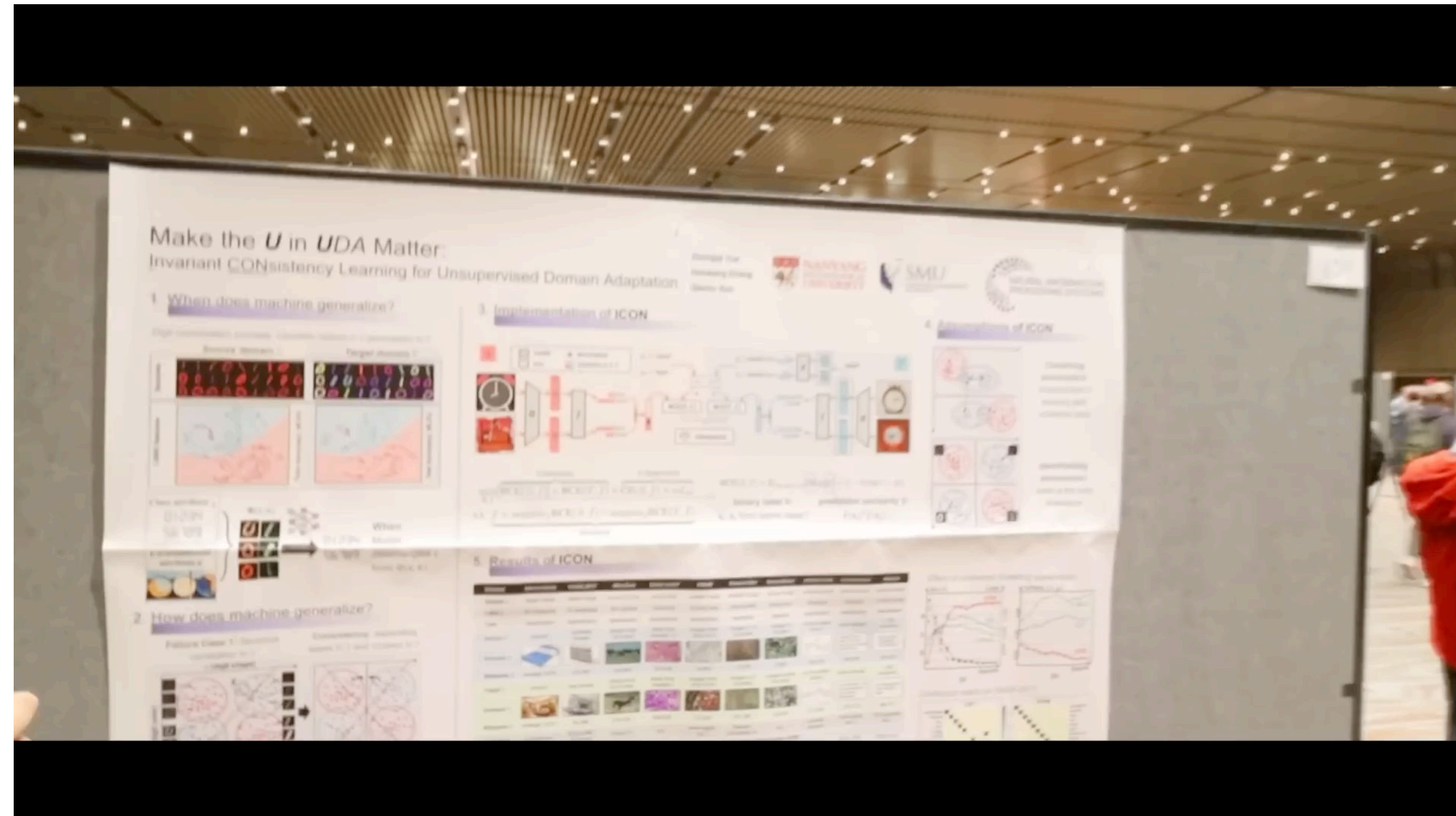
   Shorter talks

   Only 1 slide

   Direct interaction with speakers

   More flexible Q&A

In general more networking



[NeurIPS 2023 Poster Session 1 (Tuesday Evening)](https://www.youtube.com/watch?v=oUqnvQm_k9M)
https://www.youtube.com/watch?v=oUqnvQm_k9M

# POSTER SESSION: FORMALITIES PART 1

Date and time: 19.02.2025, 14:00 to 17:00

Location: OMZ U011

Poster size: A0

Plotting:

    We can plot your posters at our institute

    Please submit your posters as PDFs for printing to us no later than 17.02, 10 o'clock

All members of the group have to attend for the poster session

    If in-person is not possible, a virtual presentation can be arranged

# POSTER SESSION: FORMALITIES PART 2

Presentations should last about 15 minutes

    Afterwards 5 to 10 minutes of questions

Preliminary schedule

| Time | |
|---|---|
| 14:00 | Setting up posters |
| 14:15 | First poster session: Groups 1, 3 and 4 |
| 15:30 | Break |
| 15:45 | Second poster session: Groups 5, 6 and 7 |
| 17:00 | Estimated end |

# WHAT MAKES A GOOD POSTER?

# NOT SO GOOD POSTERS

# GOOD POSTERS

# WHAT MAKES A GOOD POSTER?

Large pictures

  Maybe even many

Little text

  But in large font

Concentrate on ONE main story