

## Methods for fast inference (§IV)

### On-device computation (§IV-A)

Model design

Model compression

Hardware

### Edge server computation (§IV-B)

Data pre-processing

Resource management

### Computing across edge devices (§IV-C)

Offloading

DNN partitioning

Edge devices plus the cloud

Distributed computing

### Private inference (§IV-D)

Add noise

Secure computation