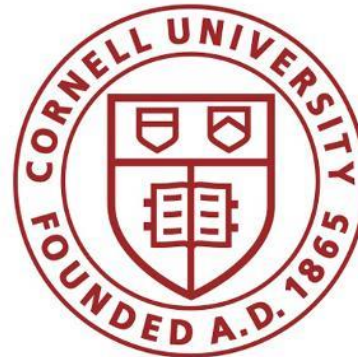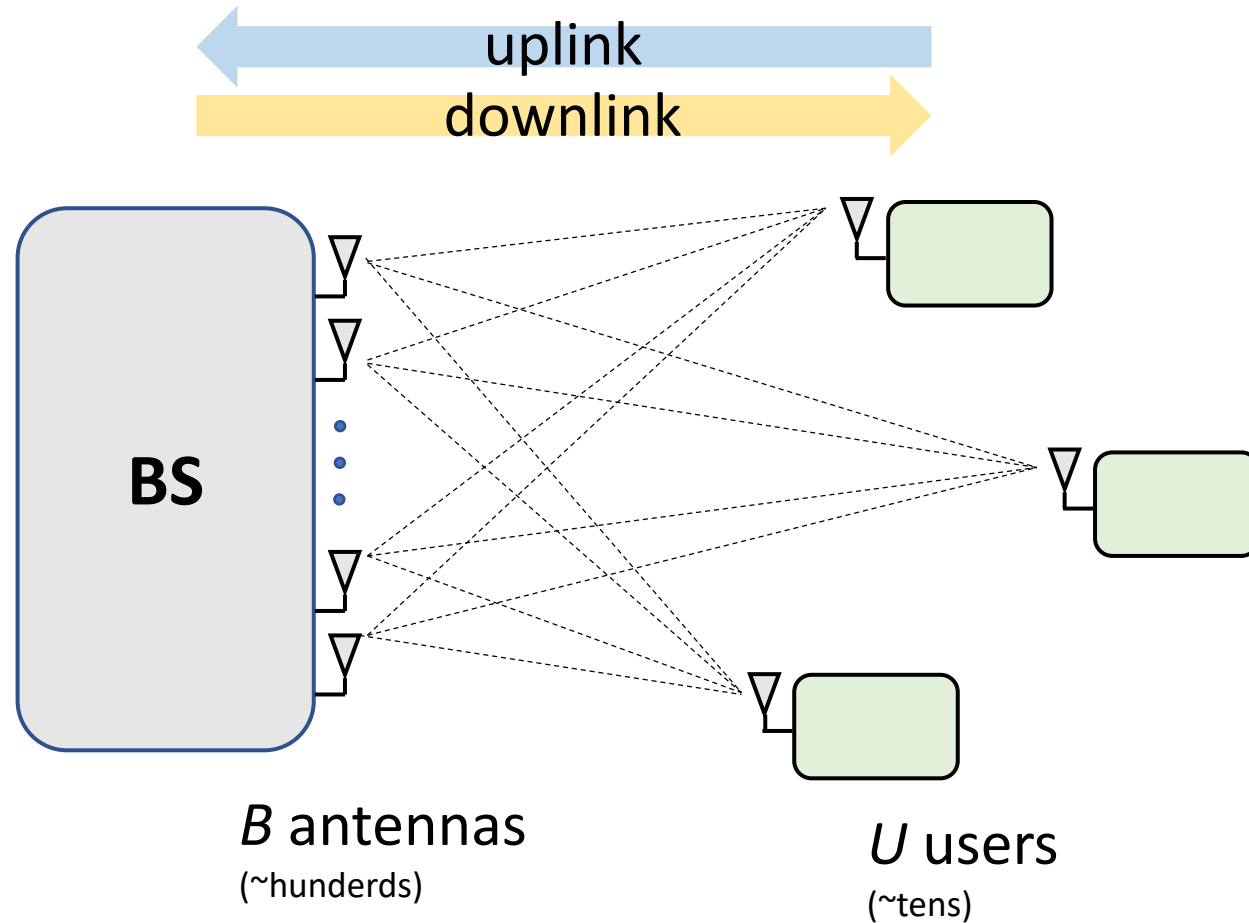# Design Trade-offs for Decentralized Baseband Processing in Massive MU-MIMO Systems

Kaipeng Li, James McNaney, Oscar Castañeda, **Chance Tarver**, Charles Jeon, Joseph Cavallaro, Christoph Studer
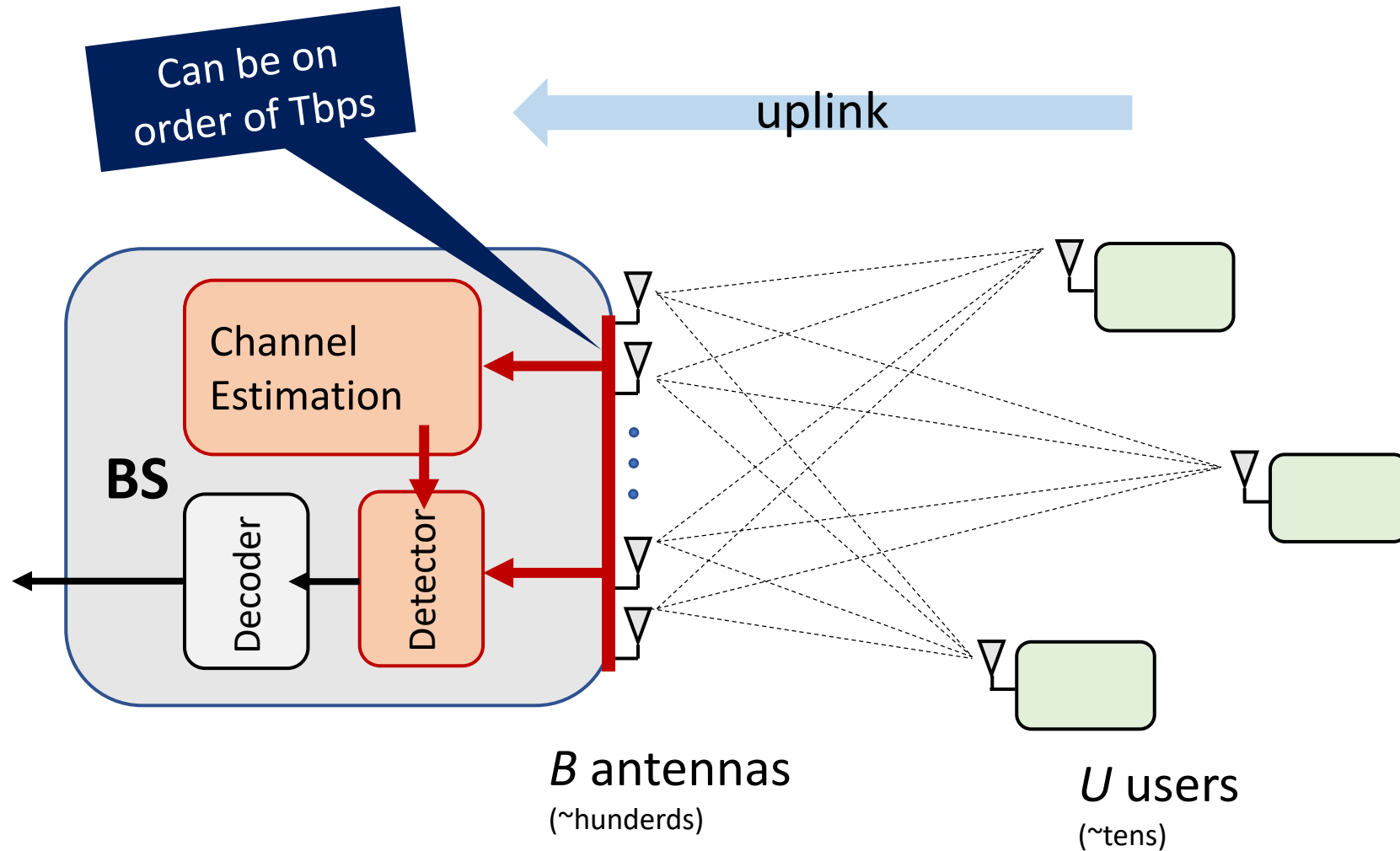
# Massive MU-MIMO systems



uplink

downlink

**BS**

*B* antennas
(~hunderds)

*U* users
(~tens)

# How do we handle this much data?



Can be on order of Tbps

uplink

**Channel Estimation**

**BS**

Decoder

Detector

*B* antennas
(~hunderds)

*U* users
(~tens)
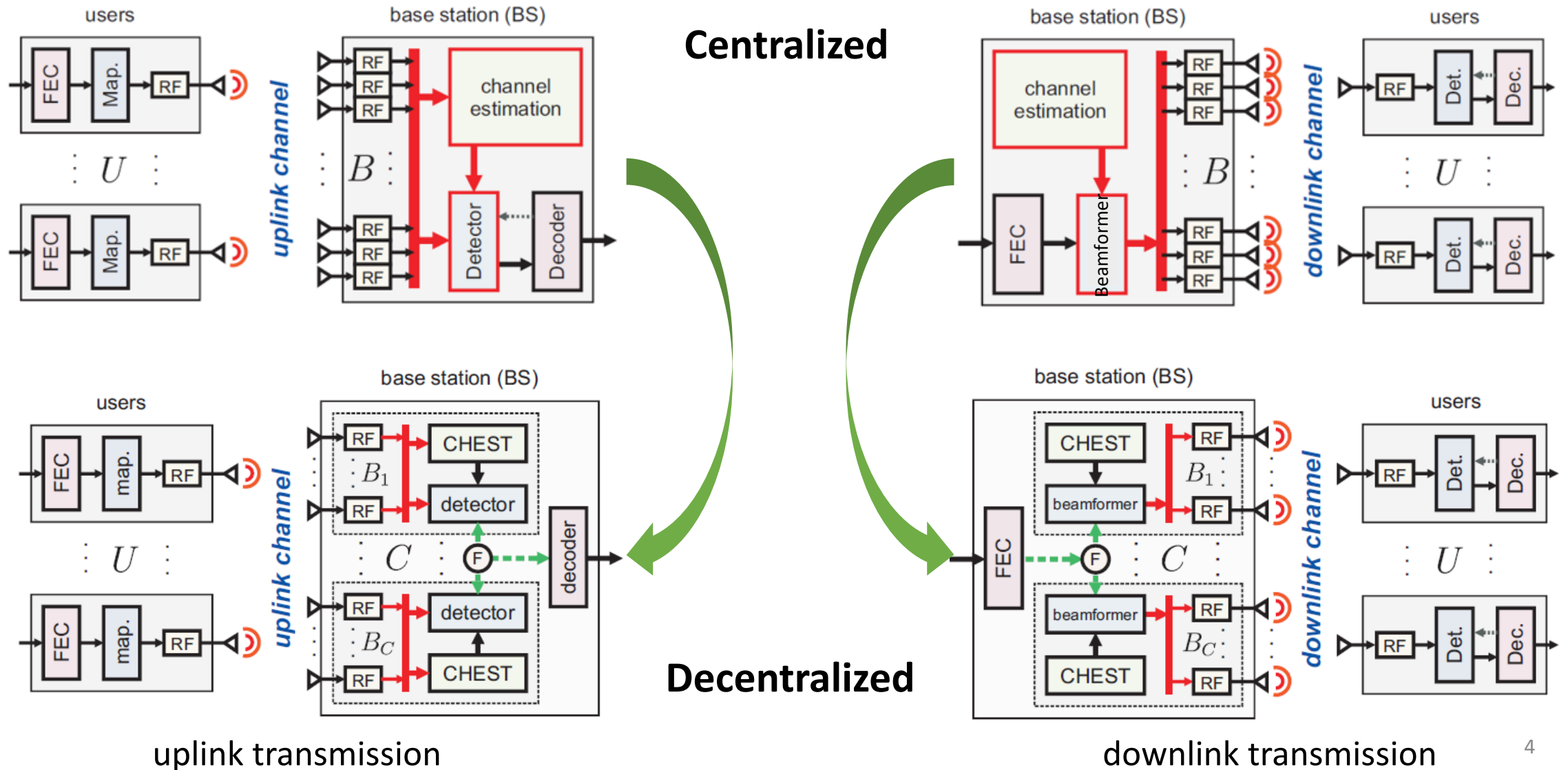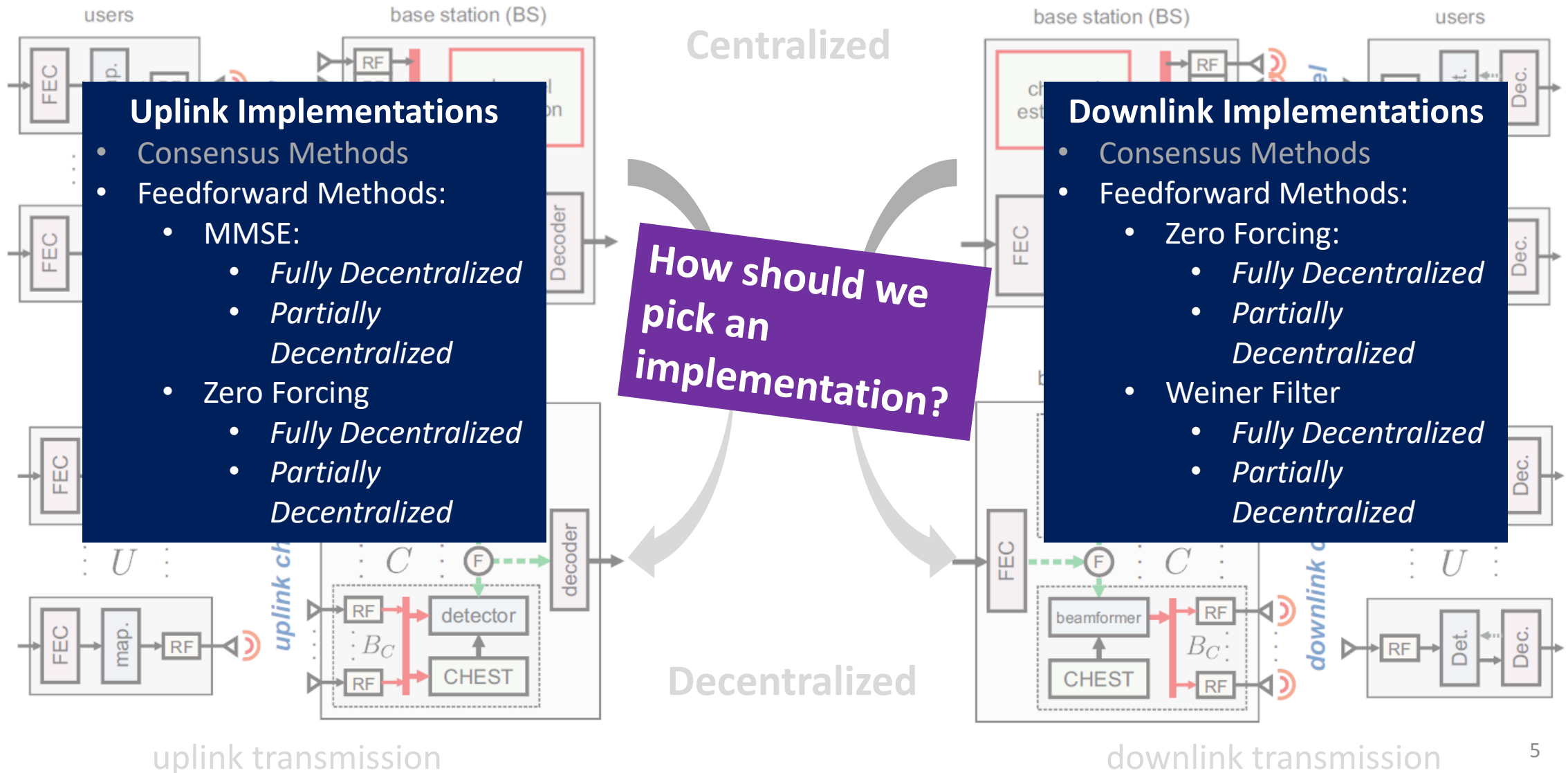
**Possible Limitations:**
- Chip I/O and interconnection bandwidth
- On-chip memory and storage
- Computing capability of modern computing fabrics

# Decentralized to resolve bottlenecks

**Centralized**



**Decentralized**

uplink transmission

downlink transmission

4

# Decentralized to resolve bottlenecks

**Centralized**

**Decentralized**

uplink transmission

downlink transmission

**Uplink Implementations**
- Consensus Methods
- Feedforward Methods:
  - MMSE:
    - *Fully Decentralized*
    - *Partially Decentralized*
  - Zero Forcing
    - *Fully Decentralized*
    - *Partially Decentralized*

**How should we pick an implementation?**

**Downlink Implementations**
- Consensus Methods
- Feedforward Methods:
  - Zero Forcing:
    - *Fully Decentralized*
    - *Partially Decentralized*
  - Weiner Filter
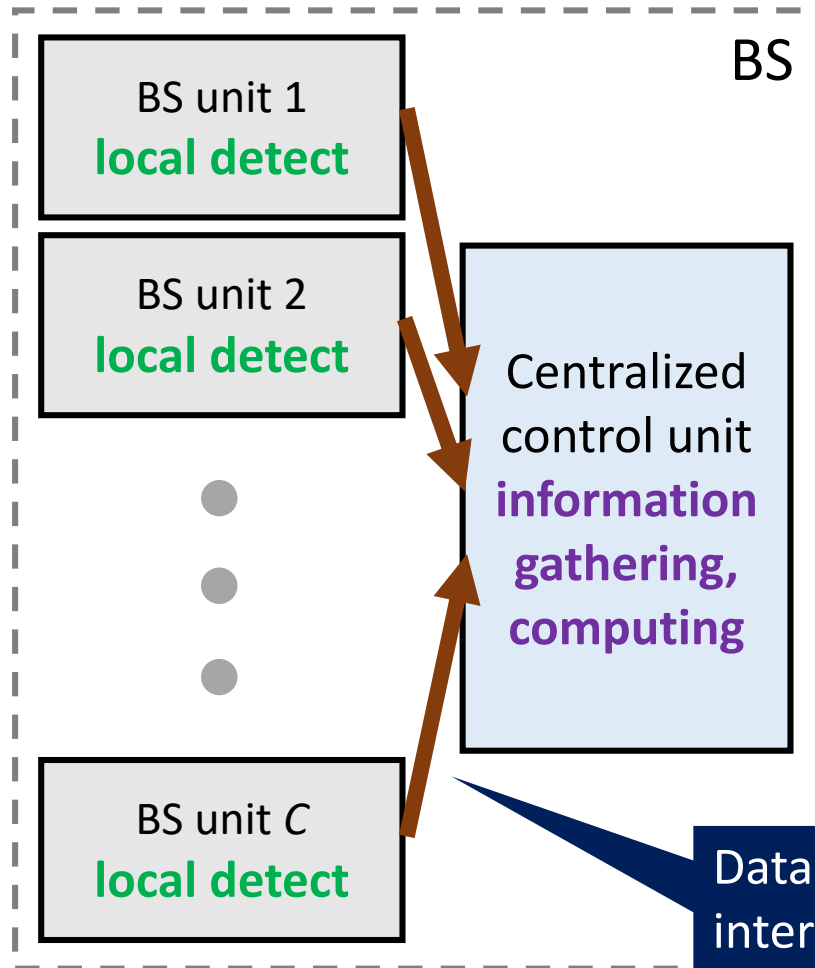    - *Fully Decentralized*
    - *Partially Decentralized*

# Outline

- Overview of decentralized architectures and algorithms
- Architecture trade-offs
- Algorithm trade-offs
- Precision trade-offs
- Practical design flow
- Conclusion

# Decentralized feedforward architecture

Feedforward local information **only once** instead of multiple rounds to centralized unit

**BS**

BS unit 1
**local detect**

BS unit 2
**local detect**

BS unit *C*
**local detect**

Centralized control unit
**information gathering, computing**

**Partially decentralized** (PD) architecture:
less local computation + more centralized computation

**Fully decentralized** (FD) architecture:
more local computation + less centralized computation

Data transfer we are interested in. Would be on PCIe, NVLink, InfiniBand.

Example: uplink system

# Uplink linear MMSE detection

Centralized linear MMSE (C-LMMSE) detection:

$$\widehat{x} = (H^H H + \frac{N_0}{E_x} I)^{-1} H^H y$$

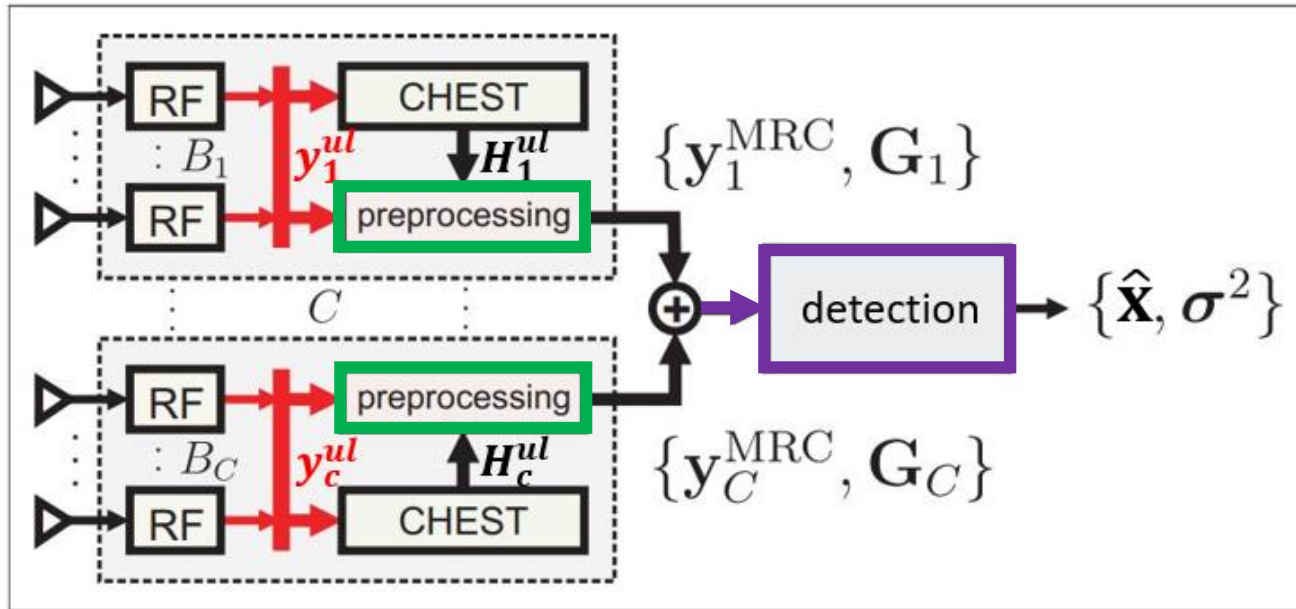$$= (\quad G \quad + \frac{N_0}{E_x} I)^{-1} y^{\mathrm{MRC}}$$

***Partially*** decentralization: ***decentralized*** matrix preprocessing + ***centralized*** detection



$\sigma^2$: error variance

# Uplink linear MMSE detection



PD-LMMSE obtains *the same* $\hat{x}$ as C-LMMSE

Complexity: $O(B_cU^2) + O(U^3)$ mults.
Data transfer size: $O(U^2)$ samples / cluster

**Centralized LMMSE (C-LMMSE):**

$$G = H^H H \qquad y^{\mathrm{MRC}} = H^H y$$

$$\hat{x} = \left(G + \frac{N_0}{E_x} I\right)^{-1} y^{\mathrm{MRC}}$$

**Partially Decentralized:**

$$G_c = H_c^H H_c \qquad y_c^{\mathrm{MRC}} = H_c^H y_c$$

$$G = \sum_{c=1}^{C} G_c \qquad y^{\mathrm{MRC}} = \sum_{c=1}^{C} y_c^{\mathrm{MRC}}$$

$$\hat{x} = \left(G + \frac{N_0}{E_x} I\right)^{-1} y^{\mathrm{MRC}}$$

# Fully decentralized (FD-) LMMSE detection



**Decentralized local detection**

$$\widehat{x}_c = (G_c + \frac{N_0}{E_x}I)^{-1}y_c^{\mathrm{MRC}}$$

**Fusion** of local $\widehat{x}_c$ using weights, $\lambda_c$ :

$$\widehat{x} = \sum_{c=1}^{C} \lambda_c \widehat{x}_c$$

Optimal $\lambda_c$ is a function of $\sigma_c$

Complexity: **O(B$_c$U$^2$) + O(U$^3$)** mults.
Data transfer size: **O(U)** samples / cluster

# Downlink Beamforming

- **Linear beamforming:**
  - Power constraint:      $\mathsf{E}[\|x\|^2] \leq \rho^2$
  - Precoding matrix:      $P$
  - Linear precoding:      $x = Ps$

- **Zero-Forcing beamforming:**
  - Precoding Matrix:    $H^H(HH^H)^{-1} = H^H \textcolor{red}{G^{-1}}$
  - Power constraint:    $\hat{x} = \rho\|\hat{x}\|_2$

- **Channel reciprocity:**
  - TDD Transmission:    $H^{dl} = (H^{ul})^H$

# Decentralized feedforward ZF beamforming



**Partially decentralized** ZF beamforming:

Set $\rho_c^2 = \rho^2/C$

$\color{green}{G_c = H_c H_c^H}$

$G = \sum_{c=1}^{C} G_c \quad z = G^{-1}s$

**Broadcast $z$ to local clusters**

$\widehat{x}_c = H_c^H z, \qquad \widehat{x}_c = \rho_c \|\widehat{x}_c\|_2$

Complexity:   $\color{green}{O(B_c U^2)} + \color{purple}{O(U^3)}$ mults.

Data transfer:  $O(U^2)$ samples / cluster

**Fully decentralized** ZF beamforming:
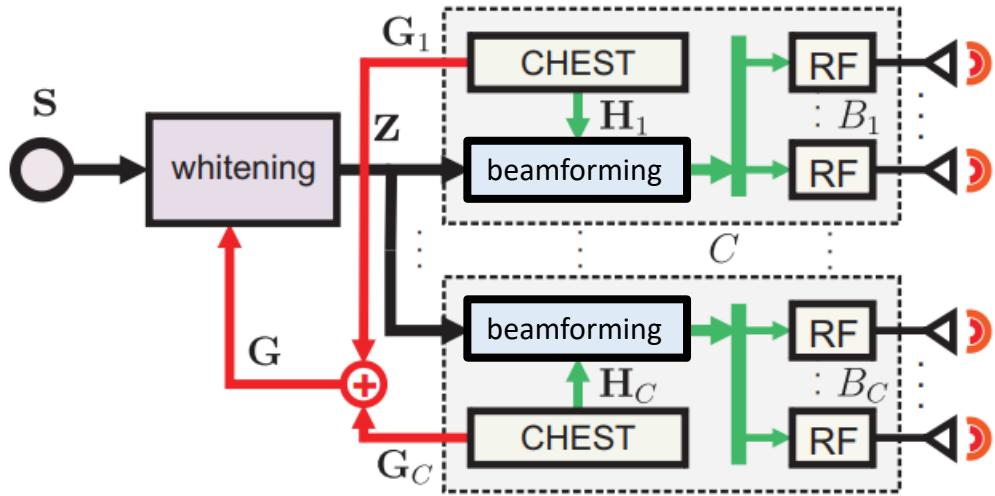
**Broadcast $s$ and set $\rho_c^2 = \rho^2/C$**

$\widehat{x}_c = H_c^H (H_c H_c^H)^{-1} s$

$\widehat{x}_c = \rho_c \|\widehat{x}_c\|_2$

Complexity:   $\color{green}{O(B_c U^2)} + \color{green}{O(U^3)}$ mults.

Data transfer:  $O(U)$ samples / cluster

# Decentralized feedforward Wiener Filter (WF) beamforming



**Partially decentralized** WF beamforming

Set $\rho_c^2 = \rho^2/C$

$$G_c = H_c H_c^H$$

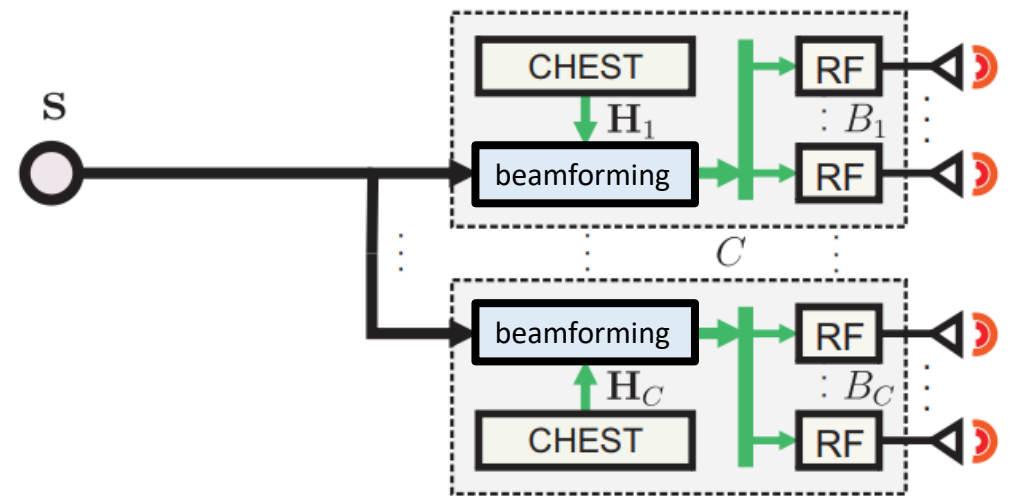$$G = \sum_{c=1}^{C} G_c \quad z = \frac{1}{\beta}(G + \gamma I_U)^{-1} s$$

**Broadcast $z$** to local BS unit

$$\widehat{x}_c = H_c^H z$$

Complexity:  $O(B_cU^2)$ + $O(U^3)$ + $O(\beta)$ mult.

Data transfer:  $O(U^2)$ samples / cluster

**Fully decentralized** WF beamforming

**Broadcast $s$** and set $\rho_c^2 = \rho^2/C$

$$P_c = \frac{1}{\beta_c} H_c^H (H_c H_c^H + \gamma I_U)^{-1} s$$

$$\widehat{x}_c = P_c s$$

Complexity:  $O(B_cU^2)$ + $O(U^3)$ + $O(\beta)$ mult.

Data transfer:  $O(U)$ samples / cluster

# Architecture Trade-offs: PD vs. FD

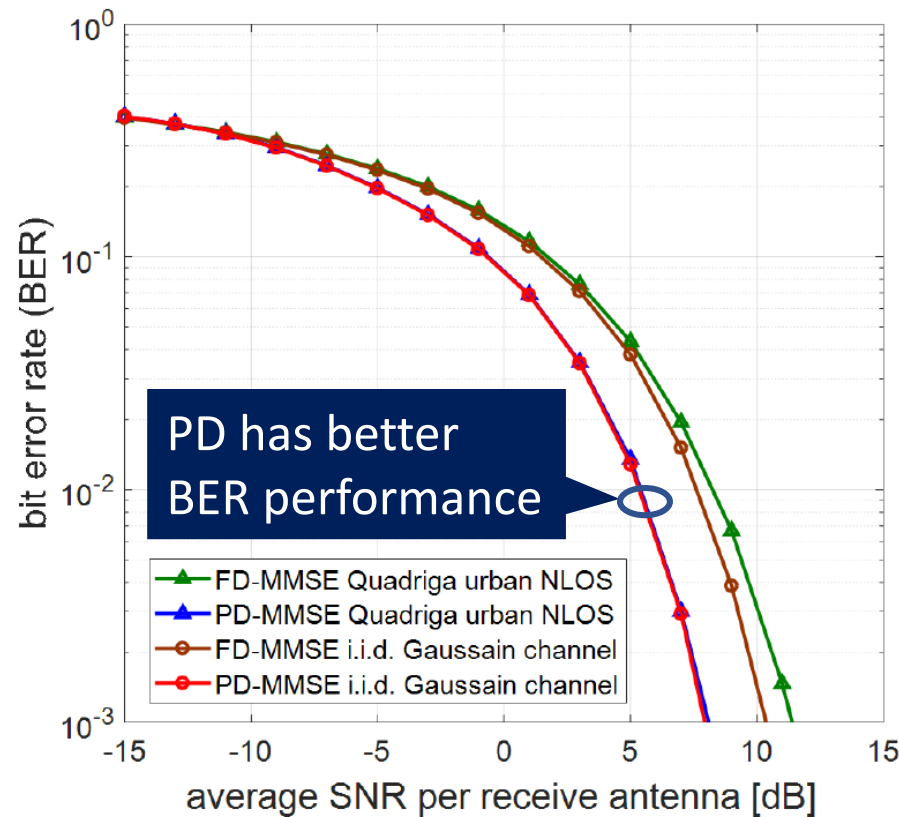**PD-MMSE and FD-MMSE: Data transfer Depends on channel coherency**

- BER: Centralized MMSE = PD-MMSE, FD-MMSE sacrifices BER

- Computation (timing) complexity: PD-MMSE = FD-MMSE

- $N_{coh}$ :  Period in which we update channel state information

$$m_{PD} = \frac{C \times (U^2 - U + 2N_{coh}U)}{N_{coh}}. \qquad m_{FD} = \frac{C \times 3N_{coh}U}{N_{coh}} = 3CU.$$
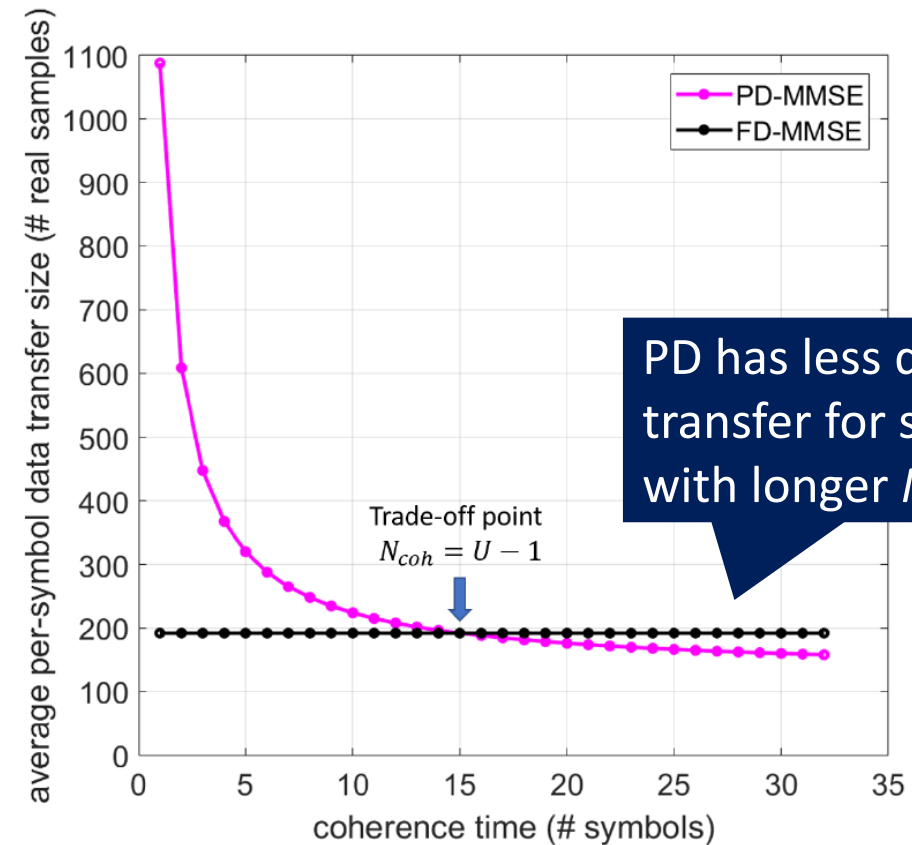
# PD vs. FD trade-off on BER and data transfer

$C=4$, $U=16$, $B_c=32$, $B=128$, 16QAM
Simple i.i.d. Gaussian channel and Quadriga NLOS urban campus channel


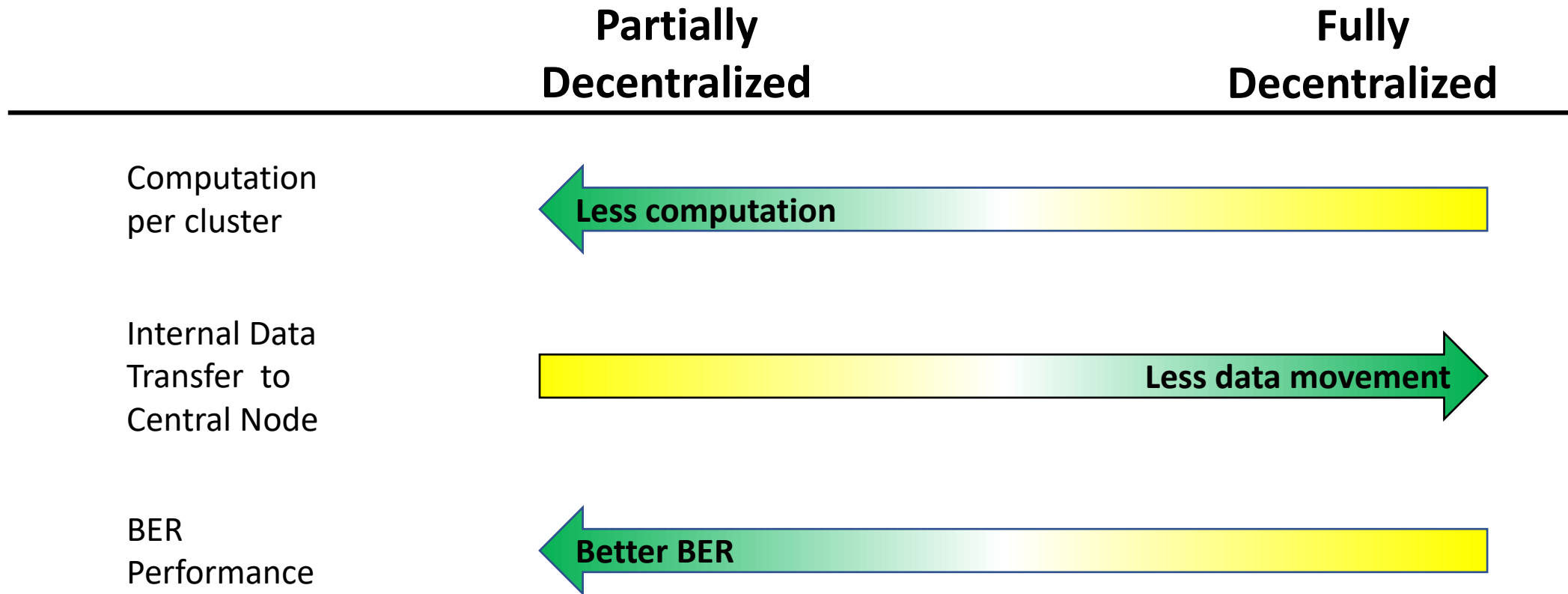
(a) BER: PD-MMSE vs. FD-MMSE

(b) Data transfer size vs. $N_{coh}$

# Decentralized feedforward architecture

Feedforward local information *only once* instead of multiple rounds to centralized unit

**Partially Decentralized**                    **Fully Decentralized**

Computation per cluster — **Less computation**

Internal Data Transfer to Central Node — **Less data movement**

BER Performance — **Better BER**

# Algorithm Trade-offs: Explicit vs. Implicit method

- Example: PD-MMSE with explicit matrix inversion vs. implicit matrix inversion

- Implicit matrix inversion $A^{-1} = (G + \frac{N_0}{E_x} I)^{-1}$ for PD-MMSE

  - $A = LL^H$ (Cholesky decomposition, $L$ is lower triangular matrix)
  - Get $z$ by solving $Lz = y^{MRC}$ using forward substitution
  - Get $\hat{x}$ by solving $L^H \hat{x} = z$ using backward substitution

- Per-symbol complexity of explicit and implicit methods depend on $N_{coh}$
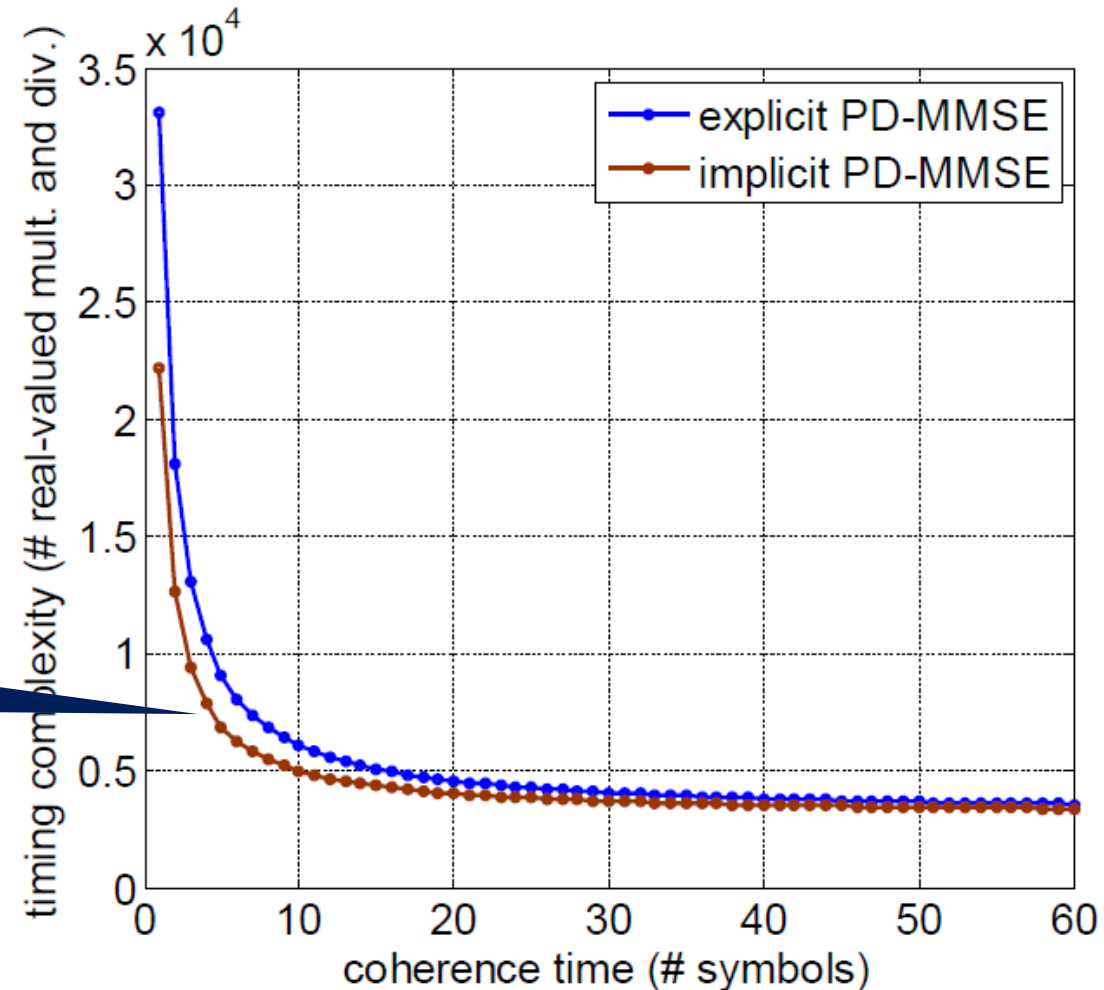
$$n_{ex} = (2B_c U^2 + \frac{10}{3} U^3 - \frac{1}{3} U)/N_{coh} + 4B_c U + 4U^2$$

$$n_{im} = (2B_c U^2 + \frac{2}{3} U^3 + \frac{1}{3} U)/N_{coh} + 4B_c U + 4U^2$$

# Complexity of explicit vs. implicit PD-MMSE

$C=4$, $U=16$, $B_c=32$, $B=128$
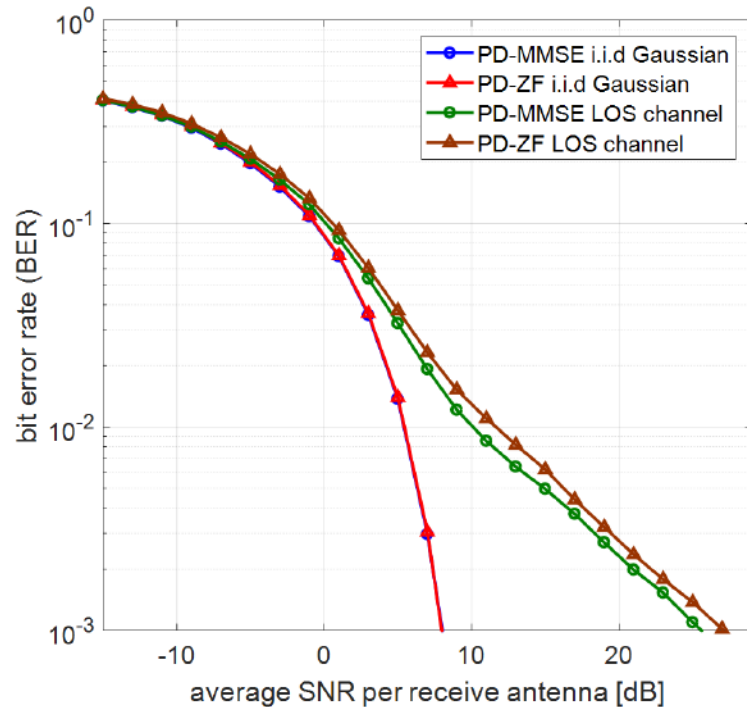
Implicit always has lower complexity
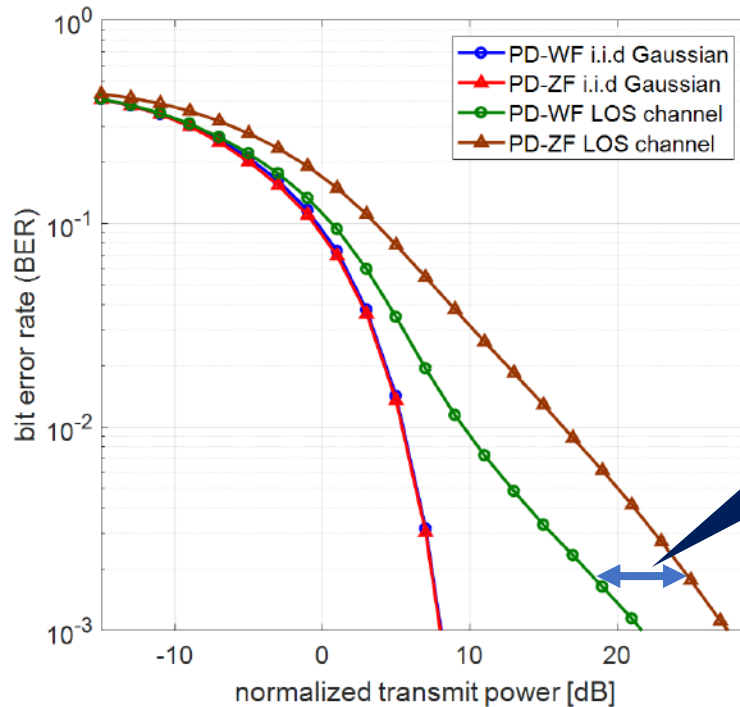
# Reusing Uplink (UL) Results for Downlink (DL)

- Channel reciprocity in TDD system: $\boldsymbol{H}^{UL} = (\boldsymbol{H}^{DL})^H$
- Gram matrix: $\boldsymbol{G}^{DL} = \boldsymbol{H}^{DL}(\boldsymbol{H}^{DL})^H = (\boldsymbol{H}^{UL})^H \boldsymbol{H}^{UL} = \boldsymbol{G}^{UL}$
- Store and reuse computed uplink results for downlink to reduce complexity
- UL MMSE detection + DL WF beamforming can only reuse $\boldsymbol{G}^{UL}$
- UL ZF detection + DL ZF beamforming can even reuse $(\boldsymbol{G}^{UL})^{-1}$

# UL and DL integration trade-offs on BER and complexity

Example: UL *PD-MMSE + DL PD-WF* integration vs. *UL ZF + DL ZF* integration
$C$=4, $U$=16, $B_c$=32, $B$=128, 16QAM, LOS channel



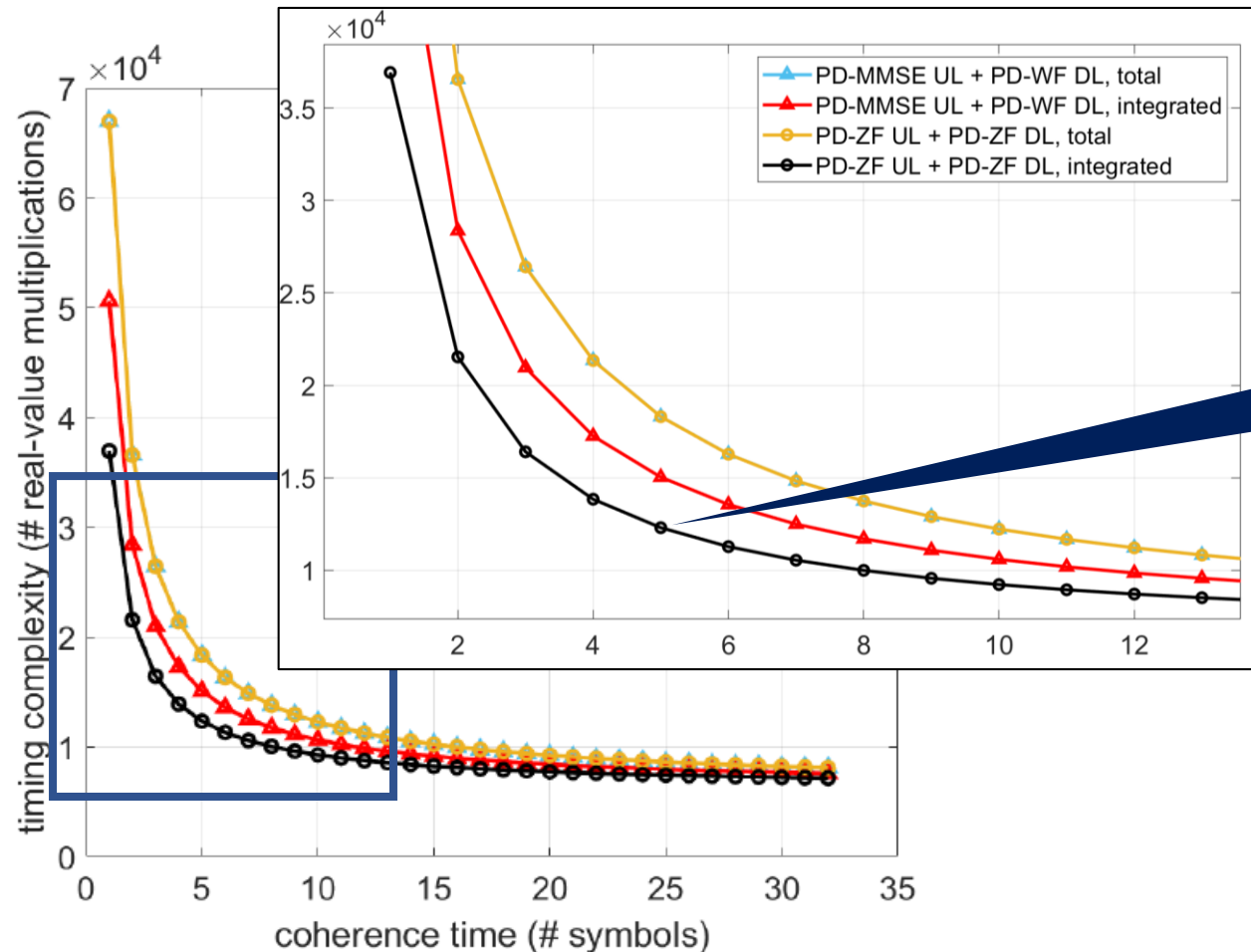(a) BER: PD-MMSE detection vs. PD-ZF detection  (b) BER: PD-WF precoding vs. PD-ZF precoding

MMSE and WF offer better performance

# UL and DL integration trade-offs on BER and complexity

Example: UL *PD-MMSE + DL PD-WF* integration vs. *UL ZF + DL ZF* integration
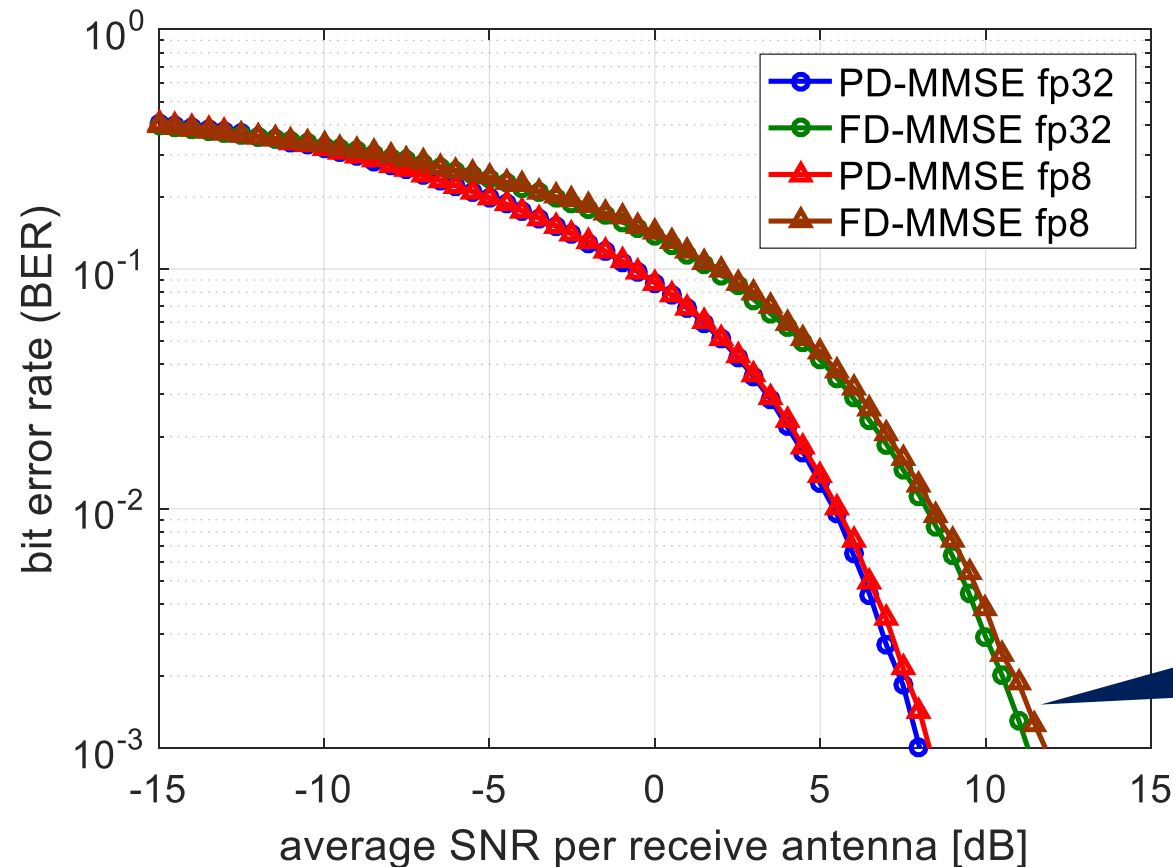$C$=4, $U$=16, $B_c$=32, $B$=128, 16QAM, LOS channel



By integrating, ZF only requires 65% of the multiplies

# Precision Trade-offs: 32-bit vs. 8bit floating point

Example: PD-MMSE and FD-MMSE

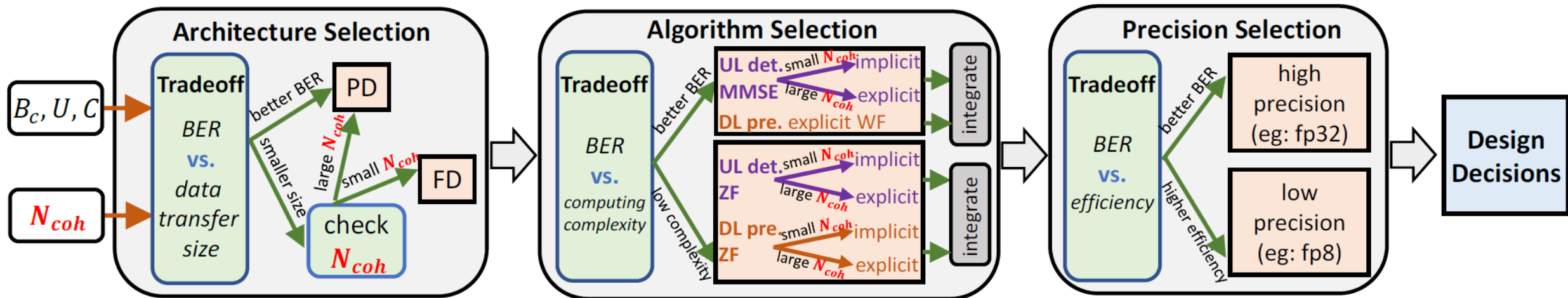C=4, U=16, $B_c$=32, B=128, 16QAM, Quadriga NLOS urban campus channel



*fp8:*  `S | 4 expo | 3 man`

FP8 is minor performance loss with significant savings

8-bit floating point reduces 4x data transfer size compared to 32-bit at only little loss of BER

# Summary of Tradeoffs

# Conclusion

- Decentralized baseband processing resolves complexity and interconnection bandwidth bottlenecks for massive MU-MIMO

- Practical massive MU-MIMO should leverage design trade-offs at different aspects:
  - Architecture trade-offs of PD and FD on BER vs. data transfer size
    - *Unless you expect very low coherence time, choose partially decentralized.*
  - Algorithms trade-offs of explicit and implicit methods on BER vs. complexity
    - *Use implicit matrix inversions whenever possible. Reuse results from uplink to downlink.*
  - Precision trade-offs of various data precision options on BER vs. efficiency
    - *Use fp16 or even fp8 unless BER is serious concern.*