

# **OEUVRE: Online Unbiased Variance-Reduced loss Estimation**

**Kanad Pardeshi, Bryan Wilder, Aarti Singh  
Carnegie Mellon University**

# Problem: Online Loss Estimation

- Consider sequential data  $\{z_t\}_{t \geq 1}$ , where  $z_t = (x_t, y_t) \sim \mathcal{D}$
- Learning algorithm  $\mathcal{A}$  trained on data, with function  $f_t : \mathcal{X} \rightarrow \mathcal{Y}$  learned using first  $(t - 1)$  samples
- **Goal:** Obtain accurate estimate of expected loss for currently learned function

$$\mathbb{E}_{(X,Y) \sim \mathcal{D}}[\ell(f_t(X), Y) \mid \mathcal{F}_{t-1}]$$

# Problem: Online Loss Estimation

- Consider sequential data  $\{z_t\}_{t \geq 1}$ , where  $z_t = (x_t, y_t) \sim \mathcal{D}$
- Learning algorithm  $\mathcal{A}$  trained on data, with function  $f_t : \mathcal{X} \rightarrow \mathcal{Y}$  learned using first  $(t - 1)$  samples
- **Goal:** Obtain accurate estimate of expected loss for currently learned function

$$\mathbb{E}_{(X,Y) \sim \mathcal{D}}[\ell(f_t(X), Y) \mid \mathcal{F}_{t-1}]$$

## Applications

Online model monitoring

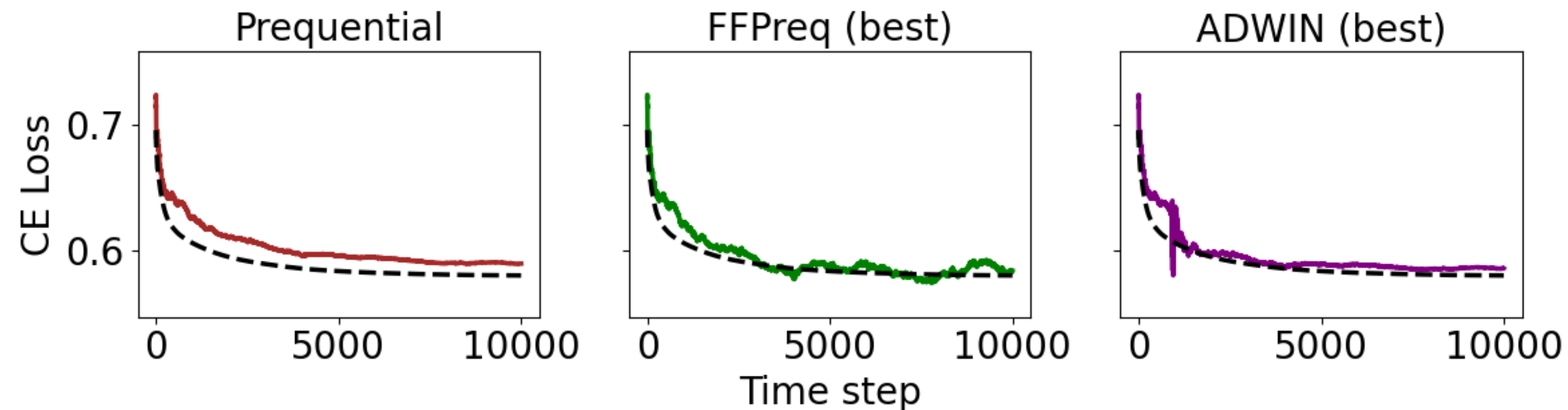
Early stopping

Online model selection

# Desirable Estimator Properties

We want an estimator with

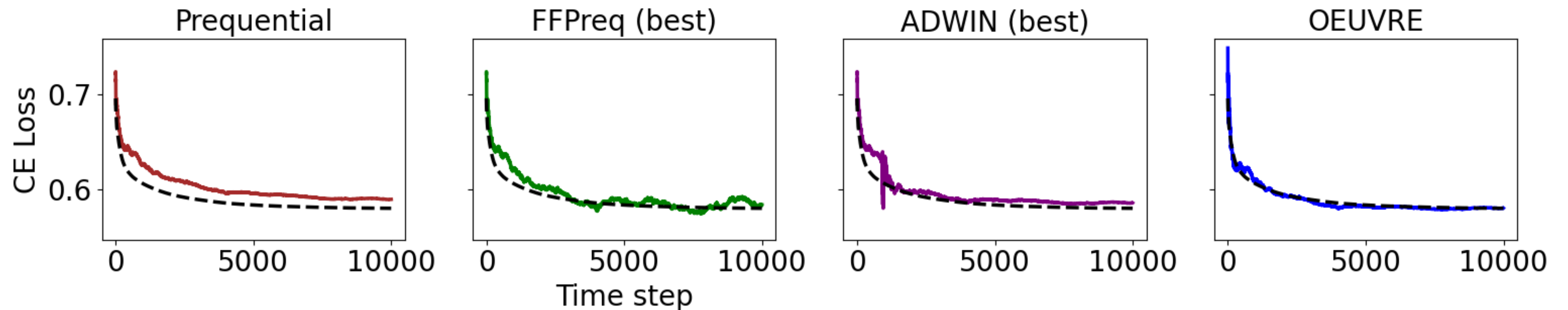
- quick and efficient updates,
- provably strong convergence properties, and
- use of properties of the learning algorithm  $\mathcal{A}$



# Desirable Estimator Properties

We want an estimator with

- quick and efficient updates,
- provably strong convergence properties, and
- use of properties of the learning algorithm  $\mathcal{A}$



Enter *OEUVRE!*

# OEUVRE Estimator

# OEUVRE Estimator

- Interleaved test-then-train: Evaluate loss of  $f_t$  on  $z_t = (x_t, y_t)$ , then train on  $z_t$

# OEUVRE Estimator

- Interleaved test-then-train: Evaluate loss of  $f_t$  on  $z_t = (x_t, y_t)$ , then train on  $z_t$
- **OEUVRE**: Evaluate loss of  $f_{t-1}$  *and*  $f_t$  on  $z_t$ , then train on  $z_t$

# OEUVRE Estimator

- Interleaved test-then-train: Evaluate loss of  $f_t$  on  $z_t = (x_t, y_t)$ , then train on  $z_t$
- **OEUVRE**: Evaluate loss of  $f_{t-1}$  *and*  $f_t$  on  $z_t$ , then train on  $z_t$ 
  - Recursive update rule

# OEUVRE Estimator

- Interleaved test-then-train: Evaluate loss of  $f_t$  on  $z_t = (x_t, y_t)$ , then train on  $z_t$
- **OEUVRE**: Evaluate loss of  $f_{t-1}$  *and*  $f_t$  on  $z_t$ , then train on  $z_t$ 
  - Recursive update rule

$$L_t = \ell(f_t(x_t), y_t) + (1 - \gamma_t) \cdot [L_{t-1} - \ell(f_{t-1}(x_t), y_t)]$$

# OEUVRE Estimator

- Interleaved test-then-train: Evaluate loss of  $f_t$  on  $z_t = (x_t, y_t)$ , then train on  $z_t$
- **OEUVRE**: Evaluate loss of  $f_{t-1}$  *and*  $f_t$  on  $z_t$ , then train on  $z_t$
- Recursive update rule

$$L_t = \ell(f_t(x_t), y_t) + (1 - \gamma_t) \cdot \boxed{L_{t-1} - \ell(f_{t-1}(x_t), y_t)}$$

Variance reduction term

# OEUVRE Estimator

- Interleaved test-then-train: Evaluate loss of  $f_t$  on  $z_t = (x_t, y_t)$ , then train on  $z_t$
- **OEUVRE**: Evaluate loss of  $f_{t-1}$  *and*  $f_t$  on  $z_t$ , then train on  $z_t$ 
  - Recursive update rule

$$L_t = \ell(f_t(x_t), y_t) + (1 - \gamma_t) \cdot [L_{t-1} - \ell(f_{t-1}(x_t), y_t)]$$

Weight term.

Set  $\gamma_t$  according to learning algorithm!

Variance reduction term

# Choosing Weight $\gamma_t$

# Choosing Weight $\gamma_t$

$$L_t = \ell_t(z_t) + (1 - \gamma_t) [L_{t-1} - \ell_{t-1}(z_t)]$$

# Choosing Weight $\gamma_t$

$$L_t = \ell_t(z_t) + (1 - \gamma_t) [L_{t-1} - \ell_{t-1}(z_t)]$$

- Choose  $\gamma_t$  to minimize an upper bound on  $\text{Var}(L_t)$

# Choosing Weight $\gamma_t$

$$L_t = \ell_t(z_t) + (1 - \gamma_t) [L_{t-1} - \ell_{t-1}(z_t)]$$

- Choose  $\gamma_t$  to minimize an upper bound on  $\text{Var}(L_t)$
- Use stability properties!

# Choosing Weight $\gamma_t$

$$L_t = \ell_t(z_t) + (1 - \gamma_t) [L_{t-1} - \ell_{t-1}(z_t)]$$

- Choose  $\gamma_t$  to minimize an upper bound on  $\text{Var}(L_t)$
- Use stability properties!

Algorithm  $\mathcal{A}$  is  $\beta_t$ -uniformly stable if

$$\|\ell(f_t(x), y) - \ell(f_{t-1}(x), y)\|_\infty \leq \beta_t$$

# Choosing Weight $\gamma_t$

$$L_t = \ell_t(z_t) + (1 - \gamma_t) [L_{t-1} - \ell_{t-1}(z_t)]$$

- Choose  $\gamma_t$  to minimize an upper bound on  $\text{Var}(L_t)$
- Use stability properties!

Algorithm  $\mathcal{A}$  is  $\beta_t$ -uniformly stable if

$$\|\ell(f_t(x), y) - \ell(f_{t-1}(x), y)\|_\infty \leq \beta_t$$

- Intuitively: bound max. change in loss

# Choosing Weight $\gamma_t$

$$L_t = \ell_t(z_t) + (1 - \gamma_t) [L_{t-1} - \ell_{t-1}(z_t)]$$

- Choose  $\gamma_t$  to minimize an upper bound on  $\text{Var}(L_t)$
- Use stability properties!

Algorithm  $\mathcal{A}$  is  $\beta_t$ -uniformly stable if

$$\|\ell(f_t(x), y) - \ell(f_{t-1}(x), y)\|_\infty \leq \beta_t$$

- Intuitively: bound max. change in loss
- Use  $\beta_t$  to find optimal  $\gamma_t$

# Choosing Weight $\gamma_t$

$$L_t = \ell_t(z_t) + (1 - \gamma_t) [L_{t-1} - \ell_{t-1}(z_t)]$$

- Choose  $\gamma_t$  to minimize an upper bound on  $\text{Var}(L_t)$
- Use stability properties!

Algorithm  $\mathcal{A}$  is  $\beta_t$ -uniformly stable if

$$\|\ell(f_t(x), y) - \ell(f_{t-1}(x), y)\|_\infty \leq \beta_t$$

- Intuitively: bound max. change in loss
- Use  $\beta_t$  to find optimal  $\gamma_t$

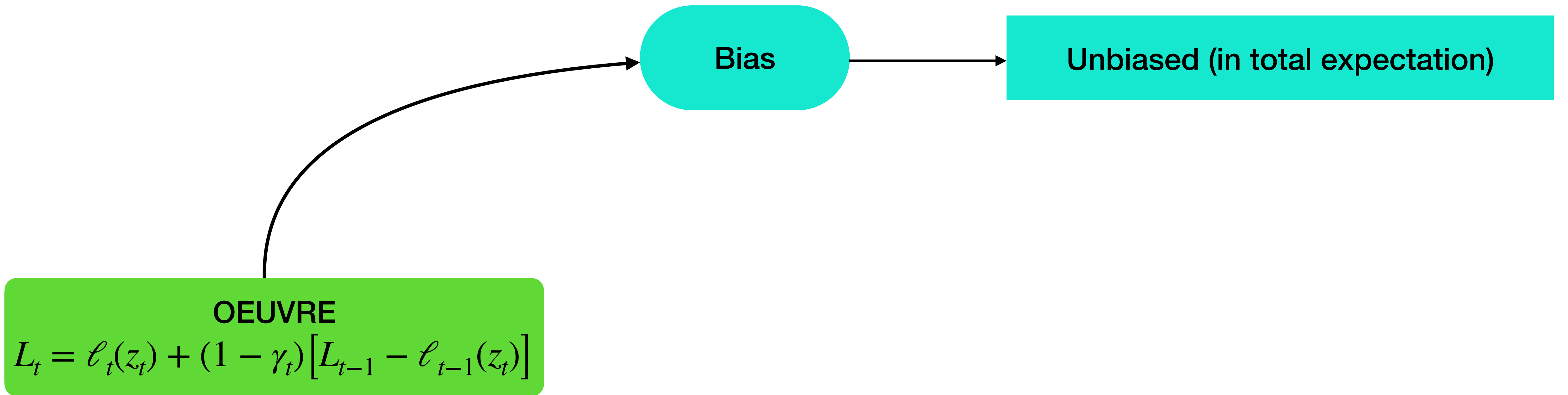
Algorithmic paradigm	Stability bound
Follow the leader (FTL)	$O(1/t)$
Follow the regularized leader (FTRL)	$O(1/\sqrt{t})$
Regularized dual averaging (RDA)	$O(1/\sqrt{t})$
Online mirror descent (OMD)	$O(\eta_t)$

# Theoretical Properties

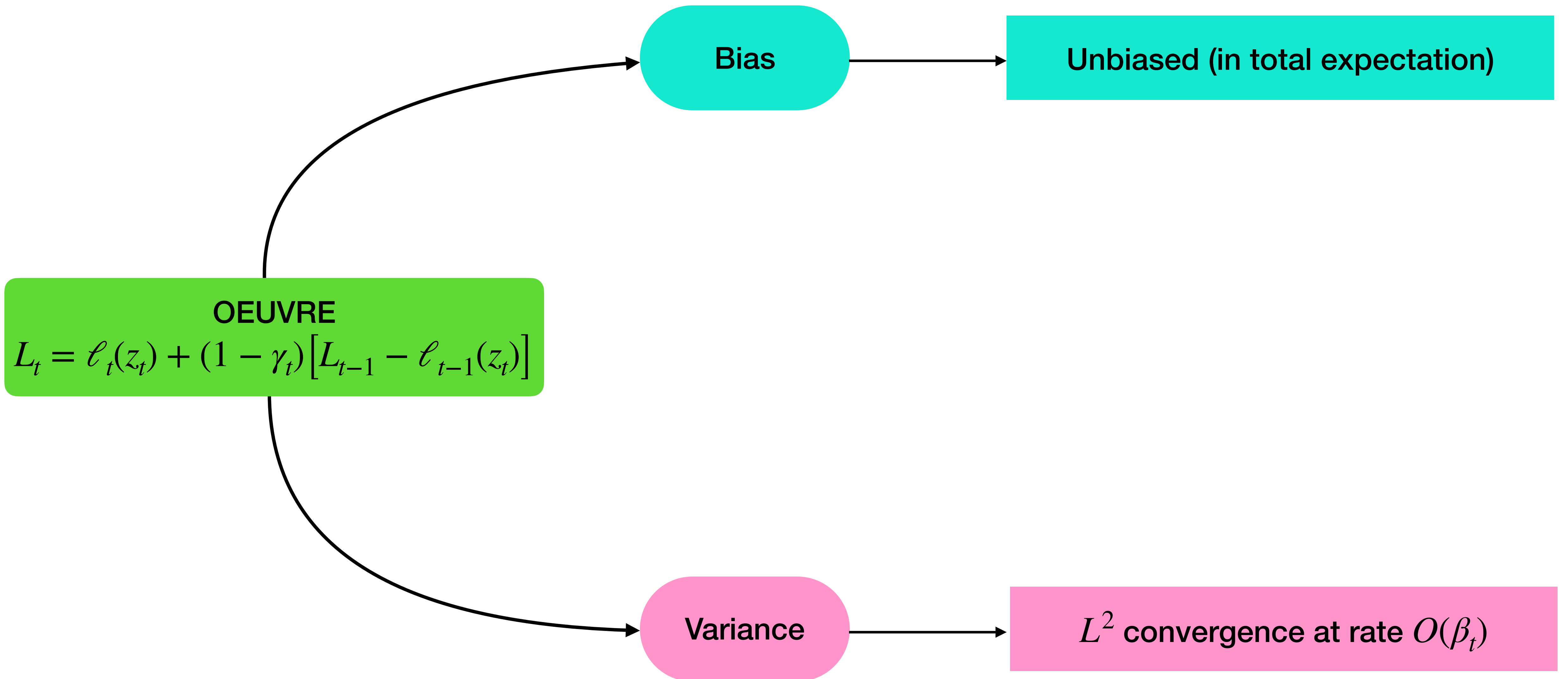
OEUVRE

$$L_t = \ell_t(z_t) + (1 - \gamma_t) [L_{t-1} - \ell_{t-1}(z_t)]$$

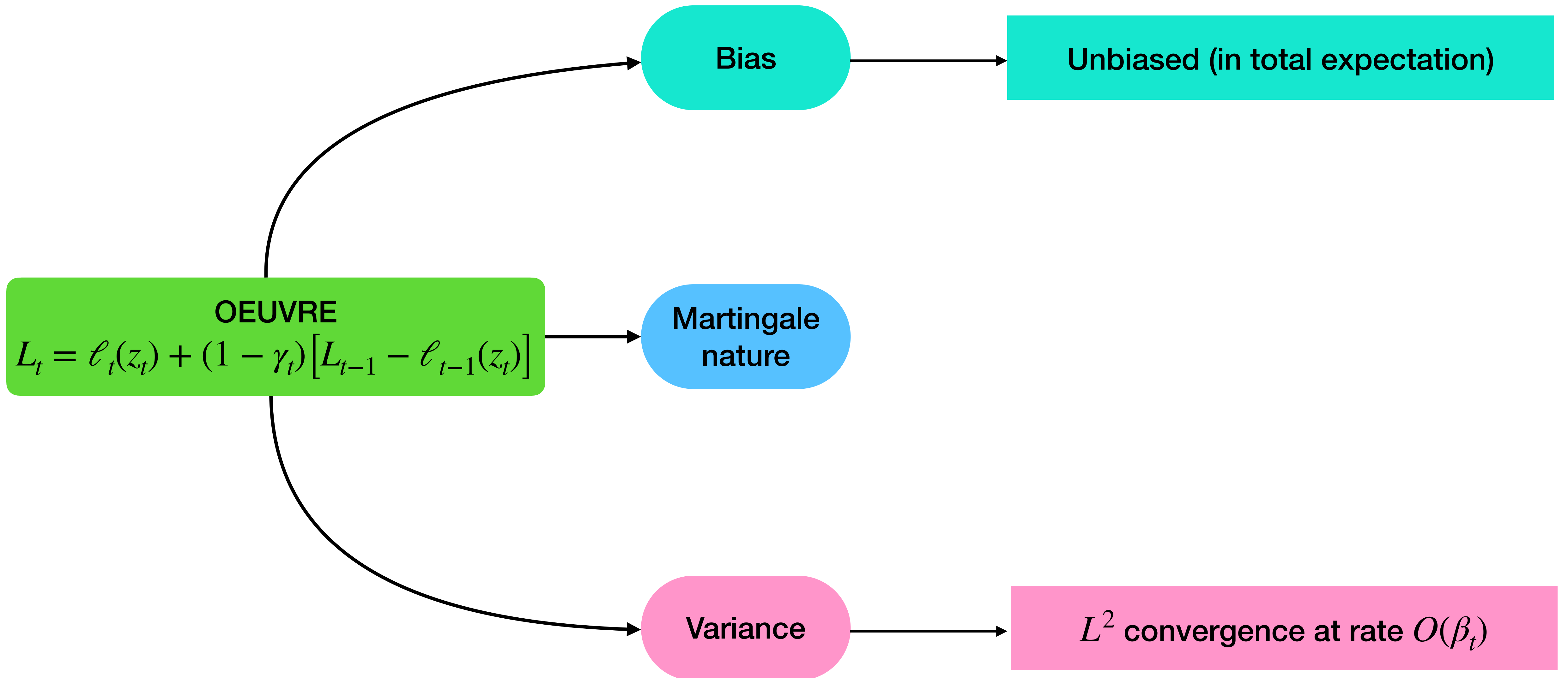
# Theoretical Properties



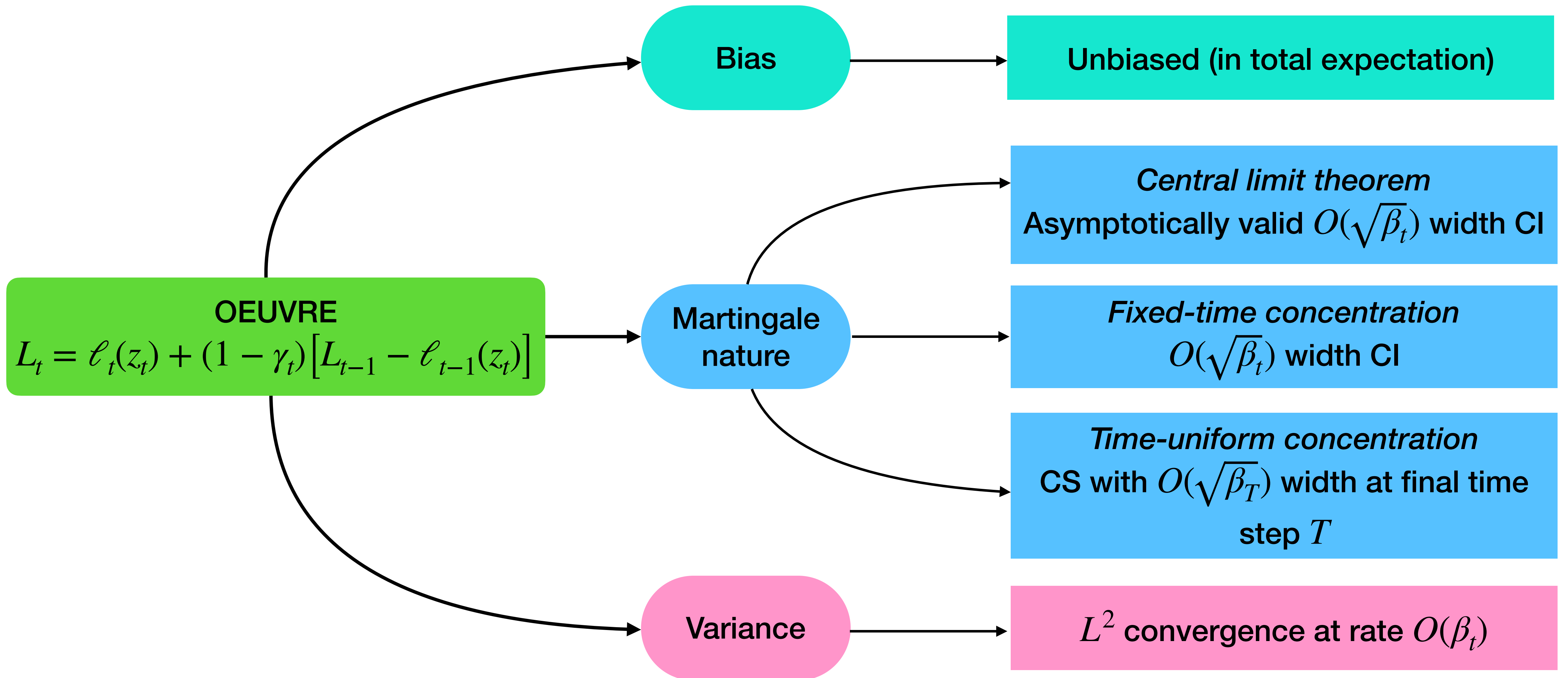
# Theoretical Properties



# Theoretical Properties



# Theoretical Properties



# Summary

- OEUVRE: extra evaluation at each time step for variance reduction
- Leverage algorithmic stability properties to design plug-and-play estimator
- Strong statistical properties (convergence, confidence intervals)
- Hyperparam-free OEUVRE: Adaptively estimate constants in stability rate
- Experiments:
  1. Linear regression
  2. Prediction with expert advice
  3. Simple neural networks

**Check out our paper!**  
Poster #173  
3:00pm

