

## Background:

- Convolutional nets are replacing recurrent nets in a growing number of sequence-to-sequence modeling applications (machine translation, audio generation, language modeling). Why?
- Both are inherently suited to modeling systems with limited long-term dependencies.
- Recurrent nets have theoretically infinite memory, but often fail to learn long sequences. Infinite memory is also usually unnecessary in practice.
- The lack of feedback elements in convolutional nets provides a computational advantage (copies of the input sequence can be processed in parallel).

## Questions:

- When do convolutional architectures provide better approximation than recurrent architectures?
- How do we quantify “limited long-term dependencies”?

## Definitions:

- i/o map* — nonlinear operator  $F : \mathcal{S} \rightarrow \mathcal{S}$  where  $\mathcal{S} := \{\mathbf{u} = (u_t)_{t \in \mathbb{Z}_+}\}$
- Right shift* —  $(R\mathbf{u})_t := u_{t-1} \mathbf{1}_{\{t \geq 1\}}$       *Window* —  $(W_{t,m}\mathbf{u})_\tau := u_\tau \mathbf{1}_{\{\max\{t-m, 0\} \leq \tau \leq t\}}$
- Causal* —  $\forall t \in \mathbb{Z}_+, (u_0, \dots, u_t) =: \mathbf{u}_{0:t} = \mathbf{v}_{0:t} \implies (F\mathbf{u})_t = (F\mathbf{v})_t$
- Time-invariant* —  $\forall k \in \mathbb{Z}_+, (FR^k\mathbf{u})_t = \begin{cases} (F\mathbf{u})_{t-k}, & \text{for } t \geq k \\ 0, & \text{for } 0 \leq t < k \end{cases}$
- Approximately finite memory* on  $\mathcal{M} \subset \mathcal{S}$  —  $\forall \epsilon > 0 \exists m \in \mathbb{Z}_+,$

$$\sup_{\mathbf{u} \in \mathcal{M}} \sup_{t \in \mathbb{Z}_+} |(F\mathbf{u})_t - (FW_{t,m}\mathbf{u})_t| \leq \epsilon$$

The smallest  $m$  such that this holds is denoted  $m_F^*(\epsilon)$ .

- Set of uniformly bounded inputs —  $\mathcal{M}(R) := \{\mathbf{u} \in \mathcal{S} : \sup_{t \in \mathbb{Z}_+} |u_t| \leq R\}$

## Main Result:

**Assumption 1:** The i/o map  $F$  has approximately finite memory on  $\mathcal{M}(R)$

**Assumption 2:** For any  $t \in \mathbb{Z}_+$ , the functional  $\tilde{F}_t : \mathbb{R}^{t+1} \rightarrow \mathbb{R}$  is uniformly continuous on  $[-R, R]^{t+1}$  with modulus and inverse modulus of continuity

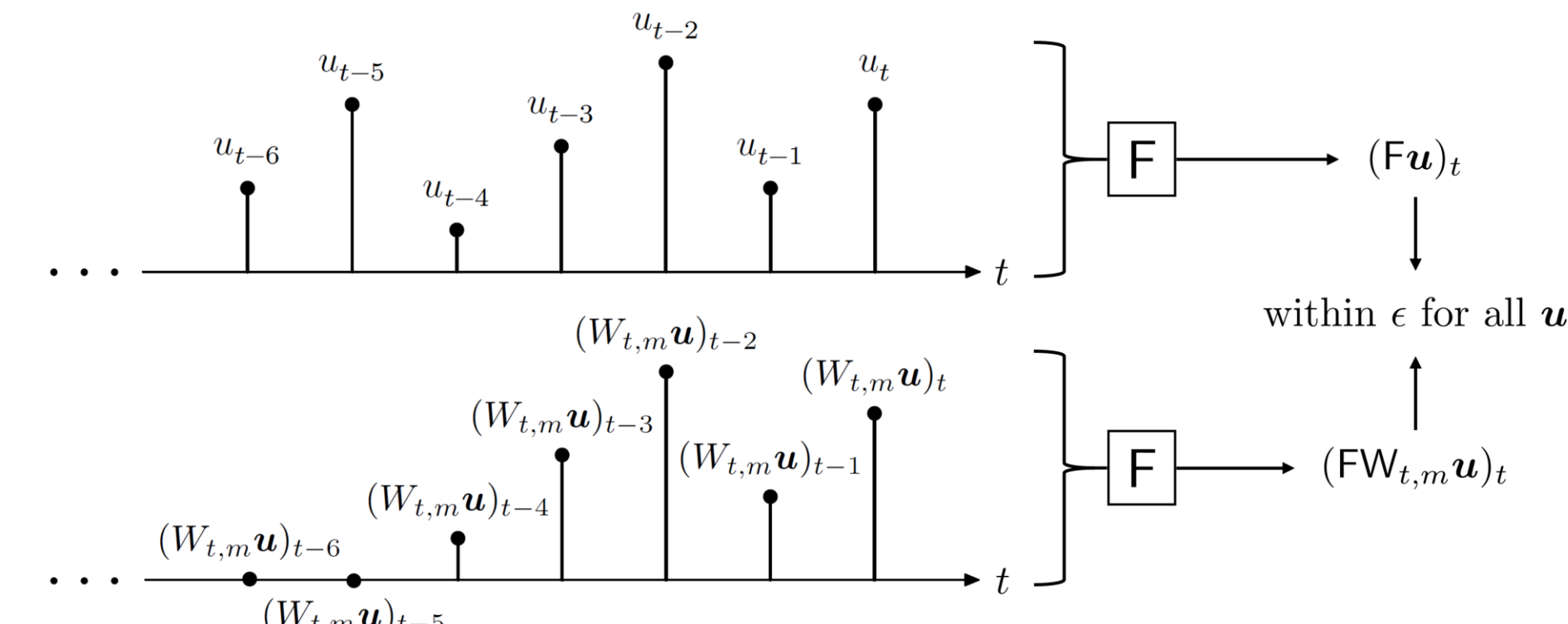
$$\omega_{t,F}(\delta) := \sup \left\{ |\tilde{F}_t(\mathbf{x}) - \tilde{F}_t(\mathbf{x}')| : \mathbf{x}, \mathbf{x}' \in [-R, R]^{t+1}, \|\mathbf{x} - \mathbf{x}'\|_\infty \leq \delta \right\}$$

$$\omega_{t,F}^{-1}(\epsilon) := \sup \left\{ \delta > 0 : \omega_{t,F}(\delta) \leq \epsilon \right\}$$

**Theorem:** Let  $F$  be a causal, time-invariant i/o map satisfying Assumptions 1 and 2. Then for any  $\epsilon > 0$ ,  $\gamma \in (0, 1)$  there exists a ReLU TCN  $\hat{F}$  with context length  $m = m_F^*(\gamma\epsilon)$ , width  $m + 2$ , depth  $\left(\frac{O(R)}{\omega_{m,F}^{-1}((1-\gamma)\epsilon)}\right)^{m+2}$  such that

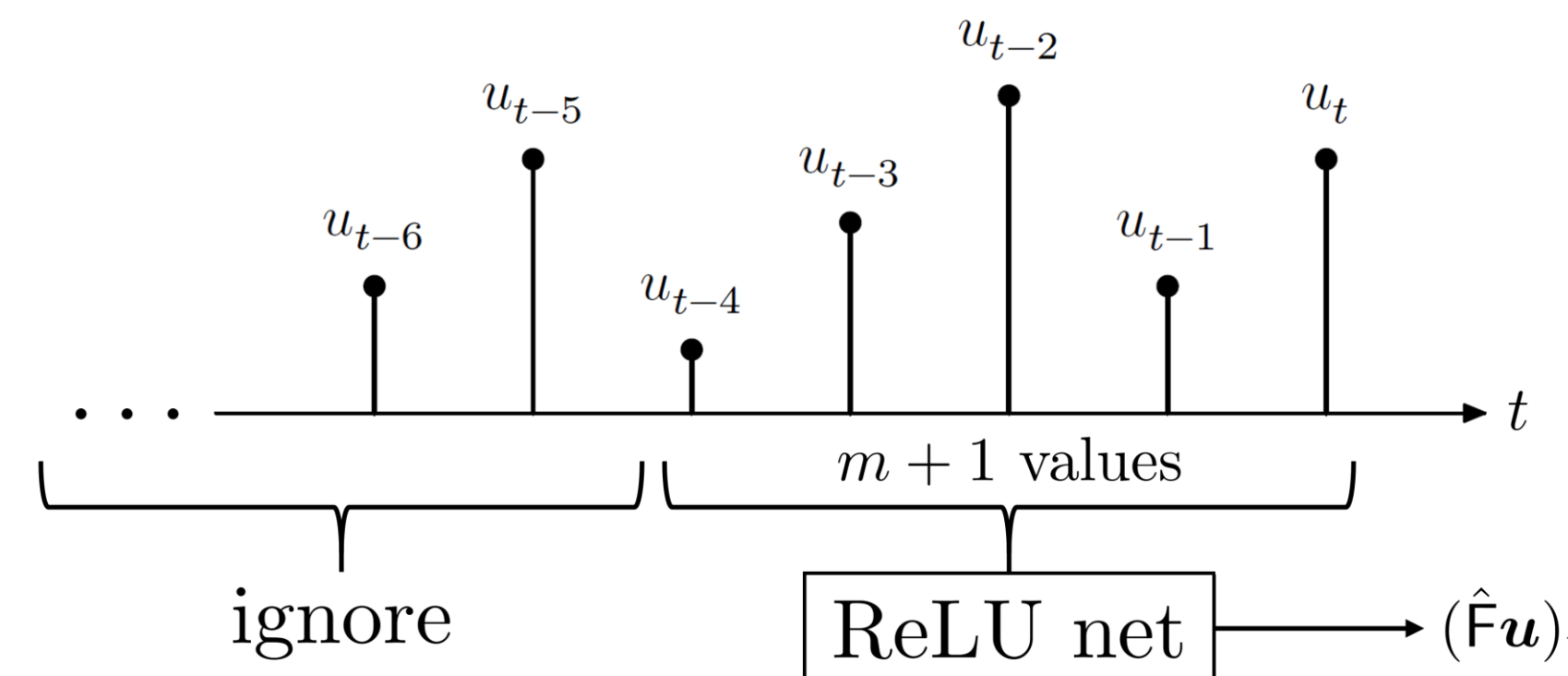
$$\sup_{\|\mathbf{u}\|_\infty \leq R} \|F\mathbf{u} - \hat{F}\mathbf{u}\|_\infty < \epsilon$$

**Remark:** The role of  $\gamma$  is to trade off context length and depth. Tuning  $\gamma$  can reduce the number of computation units needed to achieve  $\epsilon$ -accuracy.



**Figure 1 (left):**  
Definition of  
approximately  
finite memory

**Figure 2 (right):**  
Sliding context  
window of  
ReLU TCN

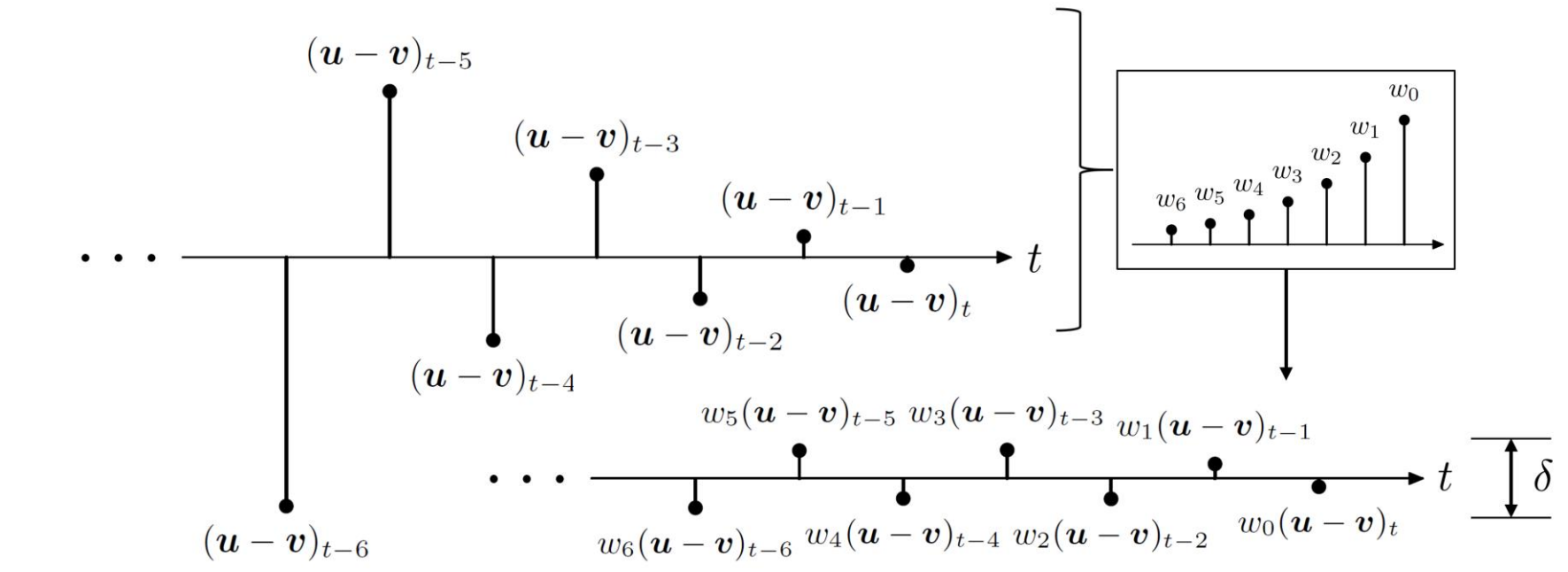


## Fading Memory:

- Set of weighting sequences —  $\mathcal{W} := \{\mathbf{w} \in \mathcal{S} : w_t \in (0, 1], w_t \downarrow 0 \text{ as } t \rightarrow \infty\}$
- $\mathbf{w}$ -fading memory on  $\mathcal{M} \subset \mathcal{S}$  —  $\forall \epsilon > 0 \exists \delta > 0 \forall \mathbf{u}, \mathbf{v} \in \mathcal{M} \forall t \in \mathbb{Z}_+$

$$\max_{s \in \{0, \dots, t\}} w_{t-s} |u_s - v_s| =: \|\mathbf{u} - \mathbf{v}\|_{\mathbf{w}} < \delta \implies |(F\mathbf{u})_t - (F\mathbf{v})_t| < \epsilon$$

**Proposition:** A continuous i/o map  $F$  satisfies Assumptions 1 and 2 if and only if it has  $\mathbf{w}$ -fading memory on  $\mathcal{M}(R)$  for arbitrary  $\mathbf{w} \in \mathcal{W}$ .



**Figure 3 (left):**  
Definition of  
fading memory

## Recurrent Systems:

Many i/o maps  $F : \mathbf{u} \mapsto \mathbf{y}$  admit state space realizations

$$x_{t+1} = f(x_t, u_t)$$

$$y_t = g(x_t)$$

where  $f : \mathbb{R}^n \times \mathbb{R} \rightarrow \mathbb{R}^n$ ,  $g : \mathbb{R}^n \rightarrow \mathbb{R}$ ,  $x_0 = \xi$ . We consider systems that are

- Uniformly asymptotically incrementally stable* on  $\mathcal{M} \subset \mathcal{S}$  — there exists a function  $\beta : \mathbb{R}_+ \times \mathbb{R}_+ \rightarrow \mathbb{R}_+$  of class  $\mathcal{KL}$  such that  $\forall \mathbf{u} \in \mathcal{M}$

$$\|\varphi_{s,t}^{\mathbf{u}}(\xi) - \varphi_{s,t}^{\mathbf{u}}(\xi')\| \leq \beta(\|\xi - \xi'\|, t - s)$$

where  $\varphi_{s,t}^{\mathbf{u}}(\xi)$  is the solution to the system above for  $x_s = \xi$  and input  $\mathbf{u}$ .

**Theorem:** Assume  $f$  and  $g$  are Lipschitz,  $\varphi_{0,t}^{\mathbf{u}}(\xi)$  remains in a compact set for all  $\mathbf{u} \in \mathcal{M}$  and all  $t \in \mathbb{Z}_+$ , and  $\beta$  is summable over its second argument. Then the i/o map for the above system satisfies Assumptions 1 and 2.

**Demidovich criterion:**  $\exists P \succ 0 \exists \mu \in (0, 1) \frac{\partial f}{\partial x}^\top P \frac{\partial f}{\partial x} - \mu P \preceq 0 \implies \beta(C, t) \propto C\mu^{t/2}$