



# **Provable Benefits of Complex Parameterizations for Structured** State Space Models

Yuval Ran-Milo, Eden Lumbroso, Edo Cohen-Karlik, Raja Giryes<sup>1</sup>, Amir Globerson<sup>1,2</sup>, Nadav Cohen<sup>1</sup>

We establish formal gaps between real and complex SSMs in terms of expressiveness and practical learnability

# **NEURAL INFORMATION** PROCESSING SYSTEMS

# **Structured State Space Models (SSMs)**

An SSM is a linear RNN (linear dynamical system)

 $\mathbf{h}_t = A \mathbf{h}_{t-1} + B x_t$  $y_t = C \mathbf{h}_t$  $A \in \mathbb{R}^{n imes n}, \; B \in \mathbb{R}^{n imes 1}, \; C \in \mathbb{R}^{1 imes n}$ 

SSMs serve as the backbone for prominent neural networks like S4, Mamba, LRU

## **Real vs. Complex Parameterizations**

SSMs can have real or complex parameterizations

### Q: To what extent do complex parameterizations benefit SSMs?

Evidence in the literature is mixed (Gu et al. 2022, Orvieto et al. 2023, Gu and Dao 2023, Ma et al. 2023)

**Conjecture:** Complex parameterizations are beneficial for continuous data, not for discrete data (Gu and Dao 2023)

# **Theoretical Result I: Separation in Expressiveness**

Denote by  $n_R$  and  $n_C$  the dimensions of real and complex SSMs

Denote by t the length of the input sequence.

A Complex SSM can precisely express any mapping realizable by a real SSM if  $n_C \ge n_R$ 

The converse is not true:

### **Theorem (informally stated)**

Many mappings realizable by a 1-dimensional complex SSM cannot be approximately expressed up to time t by a real SSM unless  $n_R \approx t$ 

## **Theoretical Result II: Separation in Practical Learnability**

What if  $n_R$  and  $n_C$  are large enough to realize a mapping up to time t?

If a mapping satisfies a **mild condition**, then a real SSM requires  $\exp(t)$  parameter values, which cannot be learned via GD

In contrast,

Any mapping can be realized by a complex SSM with parameter values linear in t

This **mild condition** is satisfied by natural mappings:

Canonical Copy

**Basic oscillatory** 



 $\mathbf{h}_t = A\mathbf{h}_{t-1} + Bx_t$  $y_t = \mathfrak{R}(C\mathbf{h}_t)$  $A \in \mathbb{C}^{n imes n}, \; B \in \mathbb{C}^{n imes 1}, \; C \in \mathbb{C}^{1 imes n}$ 

### **Theorem (informally stated)**

### Theorem (informally stated)

Random (generic)

## Experiments

Parameterization Real Complex

### **Experiments with Selectivity**

Selectivity is an SSM-based architecture yielding SotA performance.

Benefits of complex parametrization in selective SSMs are more nuanced: We see benefits on some tasks, not on others

# Future Work: Theory Accounting for Selectivity

Without Selectivity: Real SSMs struggle to express oscillations unlike complex SSMs

tasks

**Hypothesis:** Selectivity allows importing oscillations from the input

Hypothesis aligns with the conjecture of Gu and Dao as Continuous data contains only low frequencies, discrete contains all (Gu et al. 2022, Orvieto et al. 2023, Gu and Dao 2023, Ma et al. 2023)

This May fully delineate benefits of complex parameterizations for SSMs



5		
	$\supset$	ļ

ors in Theoretical Setting					
Copy (↓)	<b>Random</b> $(\downarrow)$	<b>Oscillatory</b> $(\downarrow)$			
$7.7  imes 10^{-1}$	$5.3  imes 10^{-1}$	$8.1 \times 10^{-1}$			
$1.6 imes10^{-5}$	$6.3 imes10^{-5}$	$1.6 imes10^{-4}$			

Accuracies in Real-World Setting				
Parameterization	<b>CIFAR-10</b> ( <b>†</b> )			
Real	78.27%			

Comp

 $\mathbf{89.10}\%$ 

### **Complex parameterizations for SSMs significantly** improve performance

Accuracies in different tasks				
Parameterization	Сору	Induction-Head		
Real (highest) Complex (highest)	80.17% <b>95.27%</b>	<b>98.35%</b> 97.64%		

With Selectivity: Experiments suggest real SSMs are as good as complex on some