

# Learned Benchmarks for Subseasonal Forecasting

Soukayna Mouatadid<sup>1</sup>, Paulo Orenstein<sup>2</sup>, Genevieve Flaspohler<sup>3</sup>, Miruna Opreescu<sup>4</sup>, Judah Cohen<sup>5</sup>, Franklyn Wang<sup>6</sup>, Sean Knight<sup>3</sup>, Maria Geogdzhayeva<sup>3</sup>, Sam Levang<sup>7</sup>, Ernest Fraenkel<sup>3</sup>, Lester Mackey<sup>4</sup>

1



2



3



4



Microsoft

5



6



7



# Introduction

**Subseasonal Forecasting** (3-6 weeks ahead) is a crucial pre-requisite for:

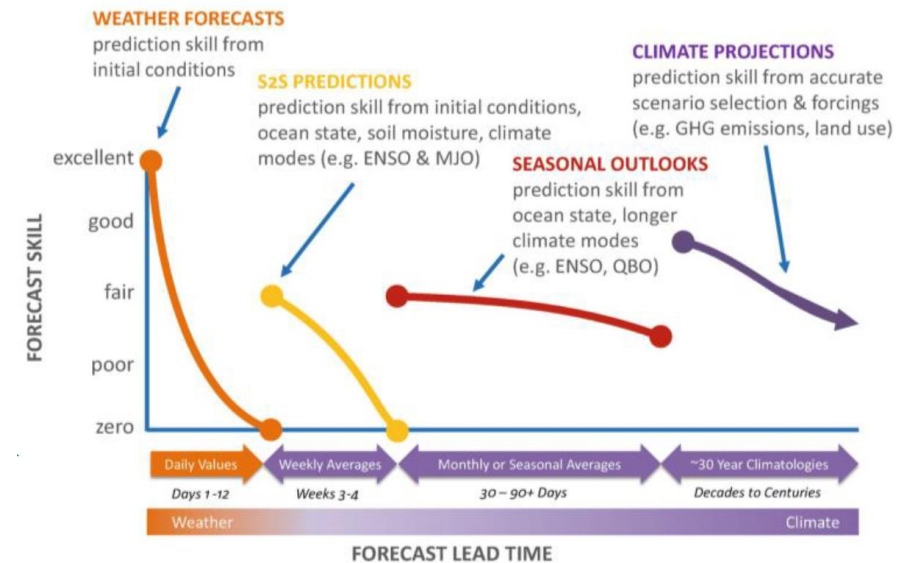
- Allocating water resources
- Preparing for droughts and floods
- Managing wildfires
- Agriculture planning

**But...**

It is a challenging forecast horizon for both meteorological and ML models

**Objective:**

- We develop a toolkit of subseasonal models that outperform operational weather models as well as state-of-the-art learning methods from the literature.



Source: Pechlivanidis and Crochemore, 2020

# Forecasting Tasks

- Target variables:
  - Average temperature (°C)
  - Accumulated precipitation (mm)
- Lead times:
  - weeks 3-4 ahead
  - weeks 5-6 ahead
- Geographical region:
  - Contiguous U.S., on a 1° x 1° grid (G = 862 grid points)

- Evaluation metrics:

- Root mean squared error:

$$\text{RMSE}(\hat{\mathbf{y}}_t, \mathbf{y}_t) = \sqrt{\frac{1}{G} \sum_{g=1}^G (\hat{y}_{t,g} - y_{t,g})^2}$$

- Skill:

$$\text{skill}(\hat{\mathbf{y}}_t, \mathbf{y}_t) = \frac{\langle \hat{\mathbf{y}}_t - \mathbf{c}_t, \mathbf{y}_t - \mathbf{c}_t \rangle}{\|\hat{\mathbf{y}}_t - \mathbf{c}_t\|_2 \cdot \|\mathbf{y}_t - \mathbf{c}_t\|_2} \in [-1, 1]$$

# Dataset

- *SubseasonalClimateUSA* dataset:
  - Regularly updated collection of ground-truth measurements and model forecasts.
  - Publicly accessible through the *subseasonal\_data* Python package
- Variables include:
  - Temperature
  - Precipitation
  - CFSv2
  - Stratospheric geopotential height
  - Madden-Julian Oscillation
  - Multivariate ENSO index
  - Pressure
  - Relative humidity
  - Sea surface temperature
  - Sea ice concentration

# Baseline Models

- Climatology

- Standard baseline for subseasonal forecasting
- Average temperature or precipitation for specific day and month over 1981-2010

- CFSv2

- Operational U.S. physics-based model from NCEP
- Main NWP baseline deployed in the U.S.

- Persistence

- Today equals tomorrow

# Learning Models

- **AutoKNN**, introduced in (Hwang et al., 2019)
- **Informer**, introduced in (Zhou, 2021)
- **LocalBoosting**, introduced in (Prokhorenkova et al., 2018)
- **MultiLLR**, introduced in (Hwang et al., 2019)
- **N-BEATS**, introduced in (Orenshkin, 2020)
- **Prophet**, introduced in (Taylor and Letham, 2018)
- **Salient 2.0**, introduced in (Schmitt, 2019)

# Our Toolkit

- **Climatology++**

- Use adaptively selected window around target day for averaging

- **CFSv2++**

- Average over range of issuance date and lead times
- Adaptively debiasing using selected window

- **Persistence++**

- Learned combination of lagged measurements with NWP

# Ensembling

- **Uniform ensemble**

- Equal-weighted average of the toolkit model forecasts
- Standard solution in the weather community

- **Online ensemble**

- Based on the AdaHedgeD online learning algorithm (Flaspohler et al., 2021)
- Results in an adaptive convex combination of base models

- **Base models**

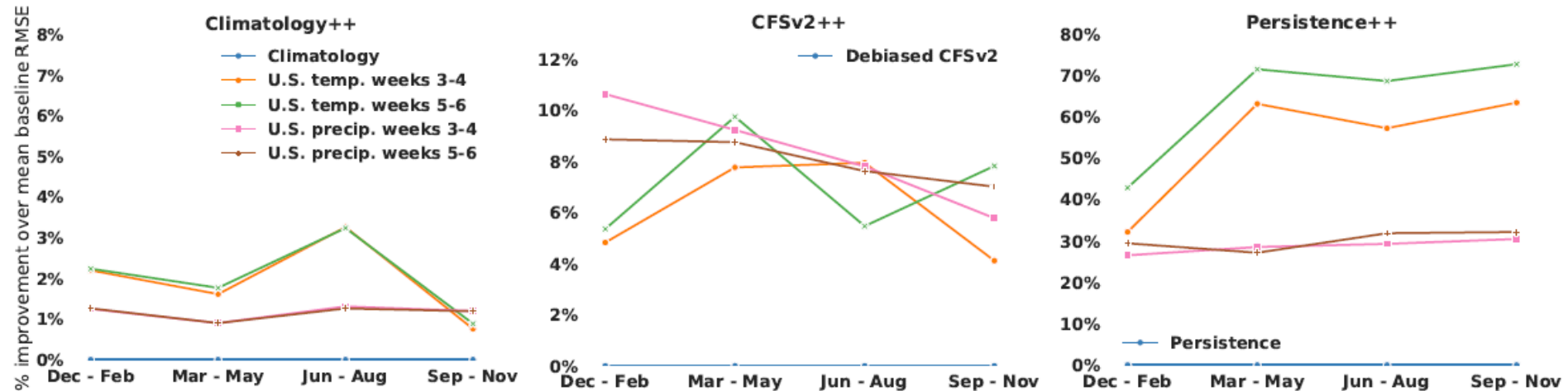
- Climatology++, CFSv2++, Persistence++, LocalBoosting, MultiLLR and Salient 2.0

# Results

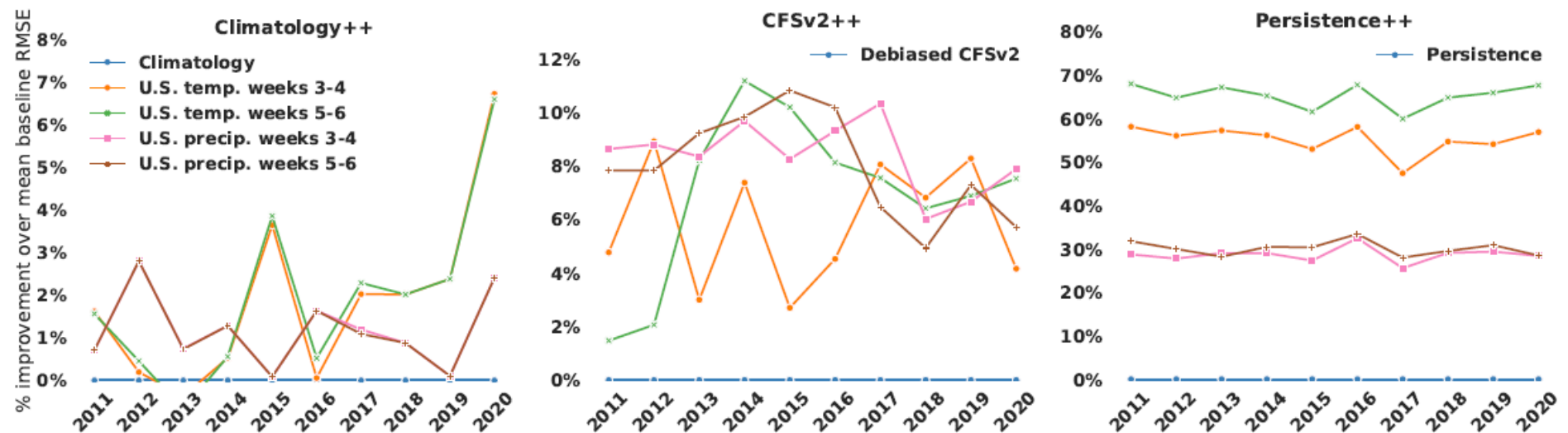
**Table 1.** Average percentage skill and percentage improvement over mean debiased CFSv2 RMSE across 2011-2020 in the contiguous U.S. The best performing model in each model group is bolded, and the best performing model overall is shown in green.

GROUP	MODEL	% IMPROVEMENT OVER MEAN DEB. CFSv2 RMSE				AVERAGE % SKILL			
		TEMPERATURE		PRECIPITATION		TEMPERATURE		PRECIPITATION	
		WEEKS 3-4	WEEKS 5-6	WEEKS 3-4	WEEKS 5-6	WEEKS 3-4	WEEKS 5-6	WEEKS 3-4	WEEKS 5-6
BASELINES	CLIMATOLOGY	0.13	<b>2.93</b>	<b>7.79</b>	<b>7.51</b>	–	–	–	–
	DEB. CFSv2	–	–	–	–	<b>23.07</b>	<b>15.98</b>	4.79	3.01
	PERSISTENCE	–109.94	–170.10	–28.27	–31.92	9.40	5.77	<b>7.84</b>	<b>7.31</b>
TOOLKIT	CLIMATOLOGY++	2.01	4.83	<b>8.86</b>	<b>8.57</b>	19.84	20.23	<b>15.44</b>	15.23
	CFSv2++	5.89	<b>7.08</b>	8.36	8.06	<b>31.40</b>	<b>27.88</b>	15.38	<b>15.29</b>
	PERSISTENCE++	<b>6.00</b>	6.43	8.61	7.89	30.19	24.91	12.39	8.88
LEARNING	AUTOKNN	0.93	3.22	7.73	7.33	12.41	9.63	5.76	5.06
	INFORMER	–39.99	–63.66	0.65	0.19	–5.17	–1.46	5.70	5.16
	LOCALBOOSTING	–0.76	–0.29	7.36	6.89	14.67	12.29	11.11	9.58
	MULTILLR	<b>2.45</b>	2.21	7.12	6.65	<b>22.37</b>	15.62	9.62	7.52
	N-BEATS	–46.71	–52.05	–19.19	–21.32	7.95	2.79	5.14	4.18
	PROPHET	1.13	<b>3.78</b>	<b>8.42</b>	<b>8.12</b>	21.13	<b>20.55</b>	<b>13.41</b>	<b>13.26</b>
	SALIENT 2.0	–6.84	–3.95	2.99	2.66	12.46	13.45	9.24	8.92
ENSEMBLES	UNIFORM TOOLKIT	6.47	7.55	9.47	<b>9.05</b>	31.96	<b>28.93</b>	18.05	<b>17.54</b>
	ONLINE TOOLKIT	<b>6.71</b>	<b>7.67</b>	<b>9.51</b>	9.04	<b>32.07</b>	28.63	<b>18.19</b>	17.30

# Results: toolkit vs baselines



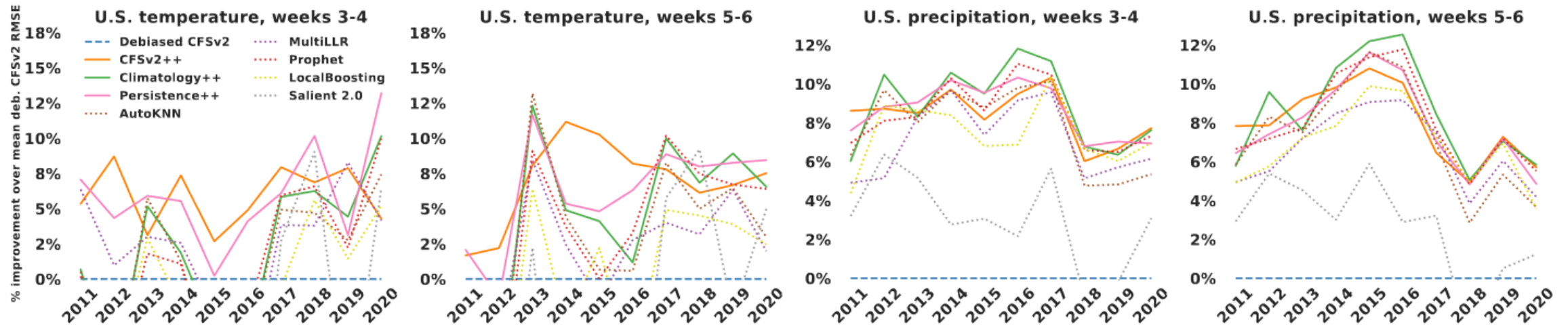
(a) RMSE improvement by quarter



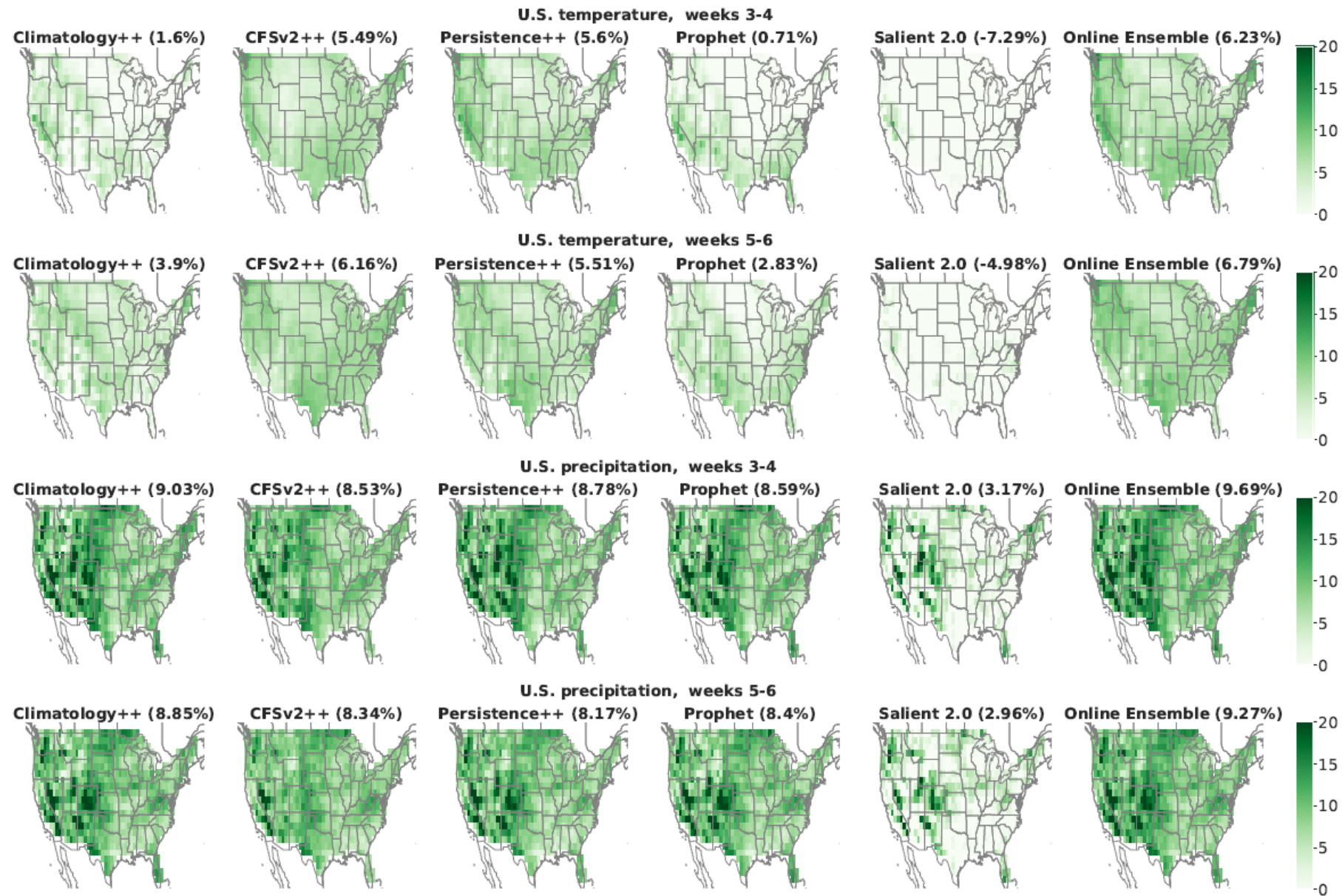
(b) RMSE improvement by year

# Results: toolkit vs learning

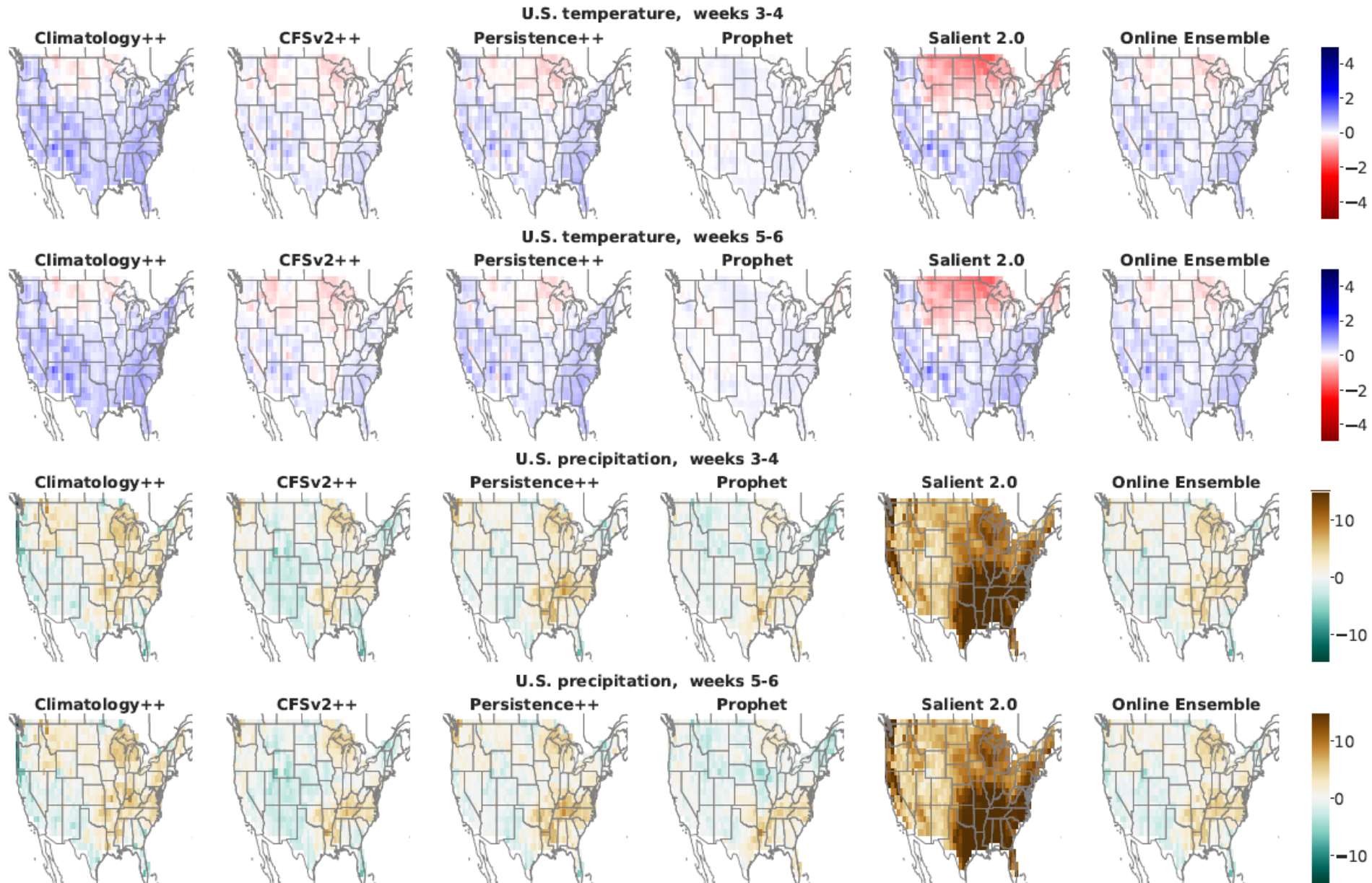
RMSE improvement by year



# Results: Percentage improvement over mean deb. CFSv2 RMSE



# Results: Mean model bias



# Results: Comparing to ECMWF

GROUP	MODEL	% IMPROVEMENT OVER MEAN DEB. CFSv2 RMSE				AVERAGE % SKILL			
		TEMPERATURE		PRECIPITATION		TEMPERATURE		PRECIPITATION	
		WEEKS 3-4	WEEKS 5-6	WEEKS 3-4	WEEKS 5-6	WEEKS 3-4	WEEKS 5-6	WEEKS 3-4	WEEKS 5-6
BASELINES	CLIMATOLOGY	<b>1.56</b>	<b>3.92</b>	<b>8.7</b>	<b>7.56</b>	–	–	–	–
	DEBIASED CFSv2	–	–	–	–	<b>22.64</b>	<b>15.71</b>	2.84	1.68
	PERSISTENCE	–105.57	–169.22	–28.05	–33.43	9.12	2.27	<b>8.11</b>	<b>6.21</b>
TOOLKIT	CLIMATOLOGY++	3.88	6.44	<b>9.79</b>	<b>8.61</b>	22.09	23.2	<b>15.34</b>	<b>15.06</b>
	CFSv2++	5.65	6.65	8.94	7.6	30.91	26.87	14.6	13.85
	PERSISTENCE++	<b>7.06</b>	<b>7.86</b>	9.06	7.57	<b>31.46</b>	<b>28.04</b>	10.03	6.61
ECMWF	DEBIASED CONTROL	–29.05	–33.25	–30.81	–31.84	18.52	13.71	0.82	3.17
	DEBIASED ENSEMBLE	<b>4.62</b>	<b>3.69</b>	<b>7.90</b>	<b>6.41</b>	<b>32.27</b>	<b>26.61</b>	<b>13.12</b>	<b>9.10</b>
ENSEMBLES	UNIFORM TOOLKIT	<b>7.43</b>	<b>8.27</b>	10.04	<b>8.77</b>	<b>32.77</b>	<b>29.75</b>	16.53	<b>15.71</b>
	ONLINE TOOLKIT	7.2	7.96	<b>10.08</b>	8.62	32.22	28.38	<b>17.19</b>	15.42

# Conclusion

- Subseasonal forecasting is a hard but fundamental problem
- Adaptive de-biasing of classical benchmarks yields sizable improvement
- Toolkit models are not only accurate, but highly scalable
- Online ensembling is highly advantageous
- Combining NWP and ML models is a powerful strategy
- All data and code are open source