



# Reanalyses for reforecast initialization

Tom Hamill

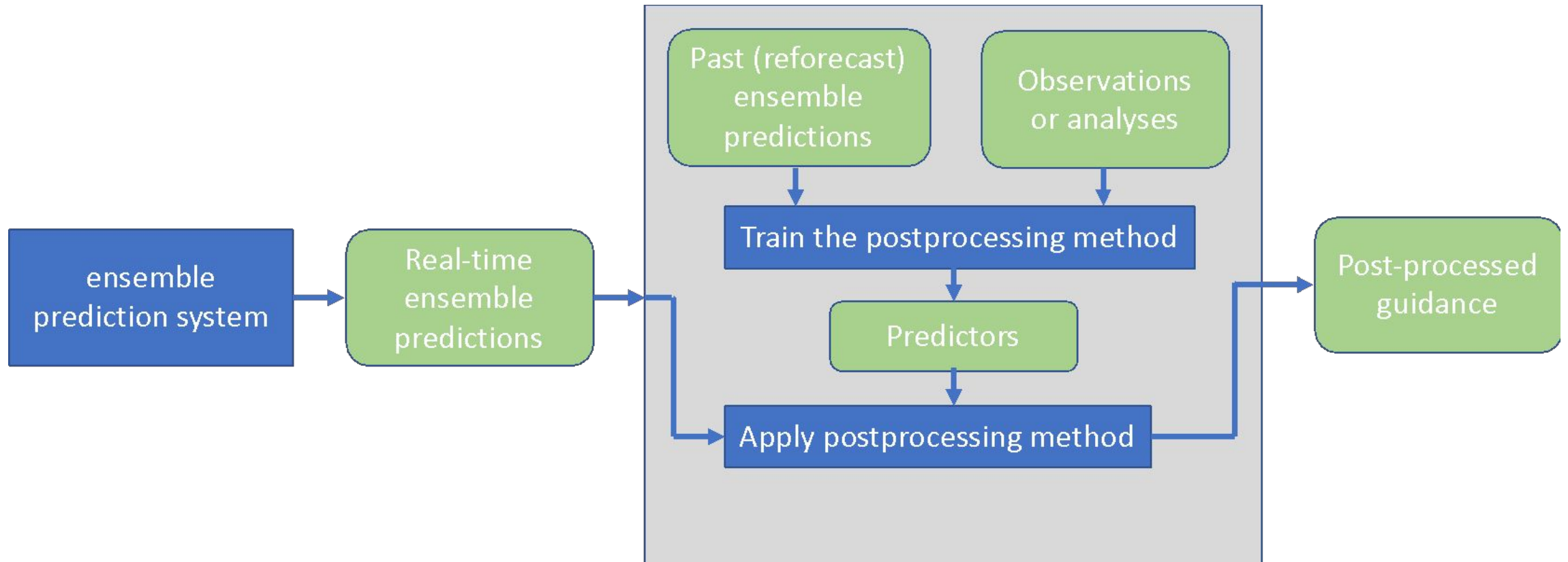
*NOAA Physical Sciences Laboratory, Boulder CO*

A presentation to US CLIVAR, 11 January 2022

# Suggested foci for the reanalysis talk from organizers (primarily address those in red).

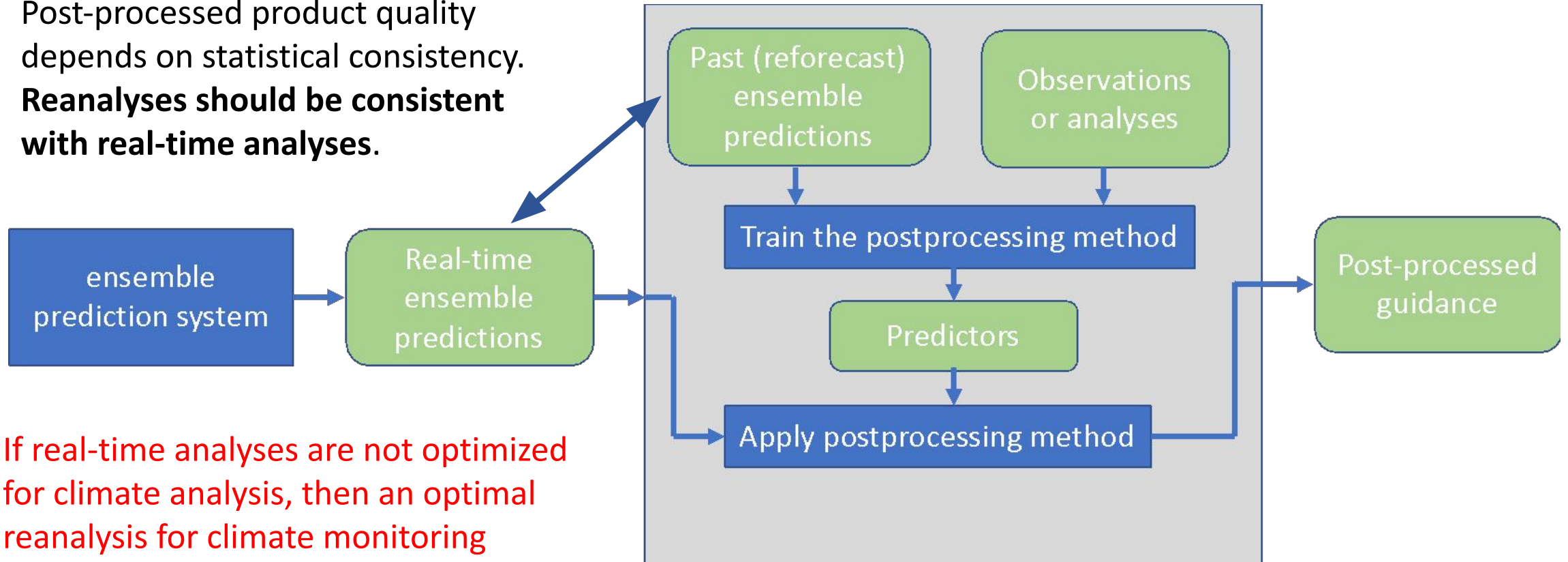
- What do you see are the most significant advances for the field of reanalysis in 5-10 years?
- What do you see are the most significant barriers to progress in the field of reanalysis?
- Which collaborations are currently working and which collaborations need to be fostered?
- What are the critical requirements for consistent Earth system reanalysis?
- What observational datasets are required to support these requirements?
- What modeling components are mature enough to enable reanalysis for your specific science question or application?
- How is uncertainty quantified for your application? Are there significant barriers for quantifying uncertainty in your field?

# Reanalysis and reforecasting as part of an integrated Unified Forecast System.



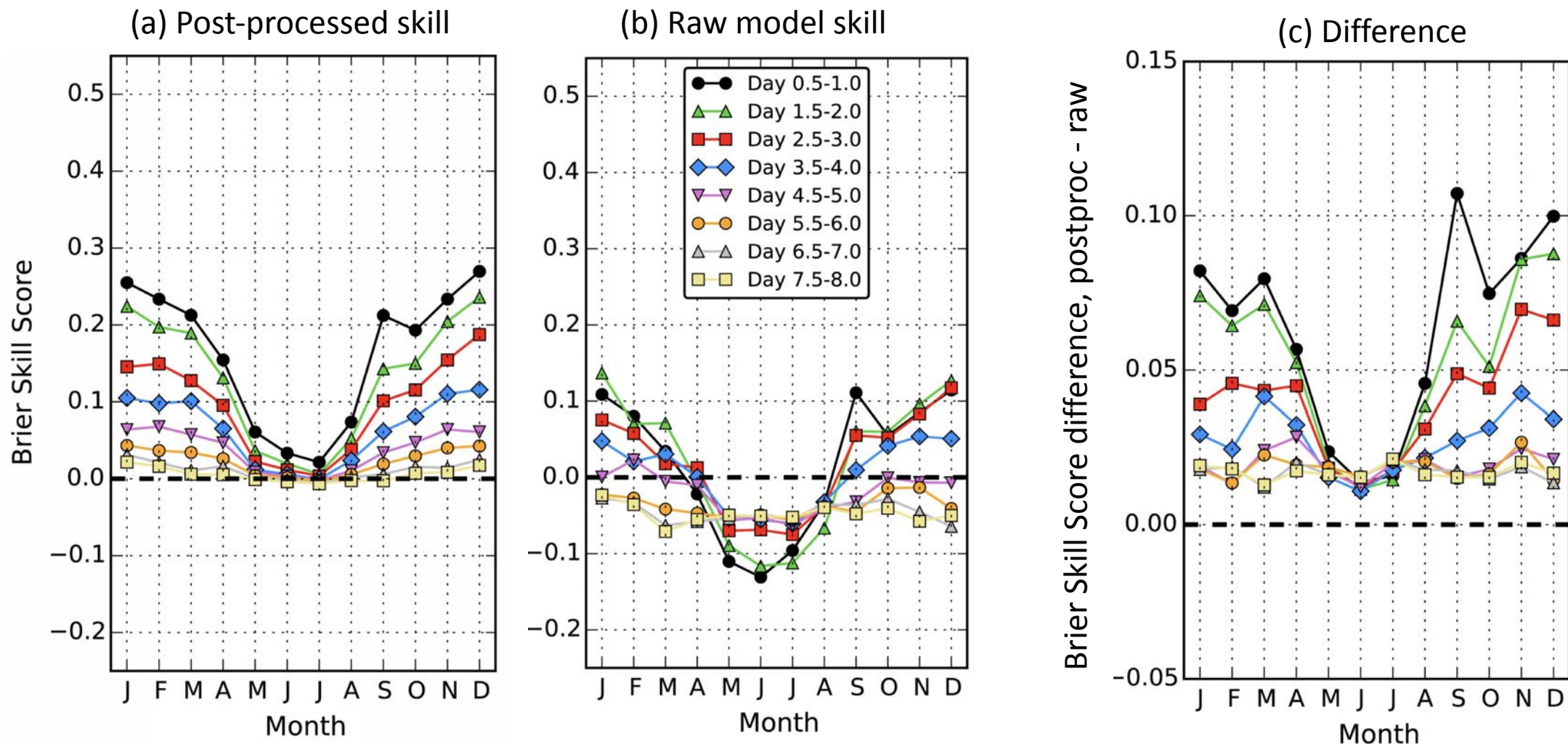
# Reanalysis and reforecasting as part of an integrated Unified Forecast System.

Post-processed product quality depends on statistical consistency. **Reanalyses should be consistent with real-time analyses.**

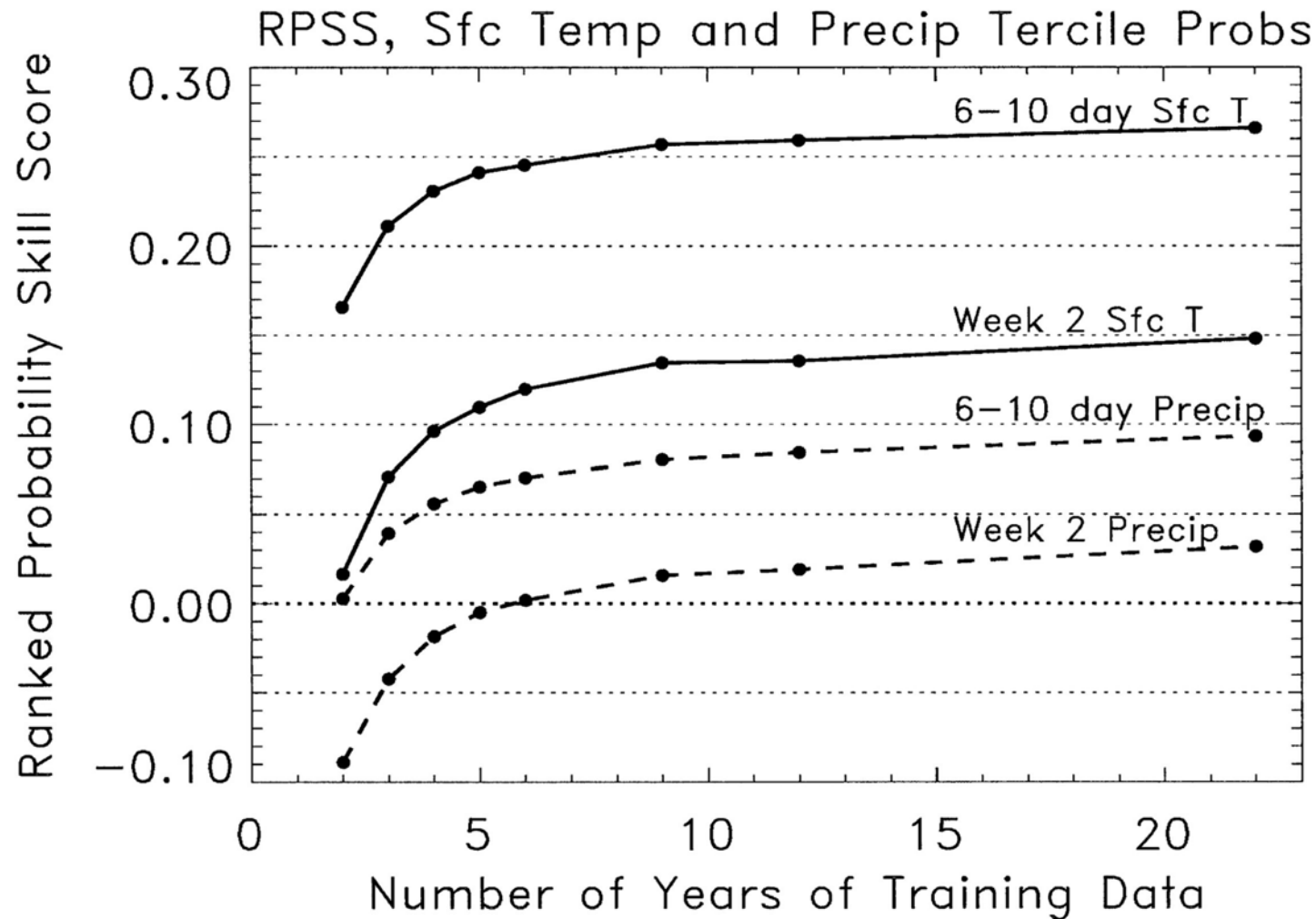


If real-time analyses are not optimized for climate analysis, then an optimal reanalysis for climate monitoring may be different from reanalysis for reforecast initialization.

# Impact of long training data set on reforecasts postprocessing. Heavy precipitation skill (>25 mm/12h)



# Impact of length of training data set (first-generation reforecasts)



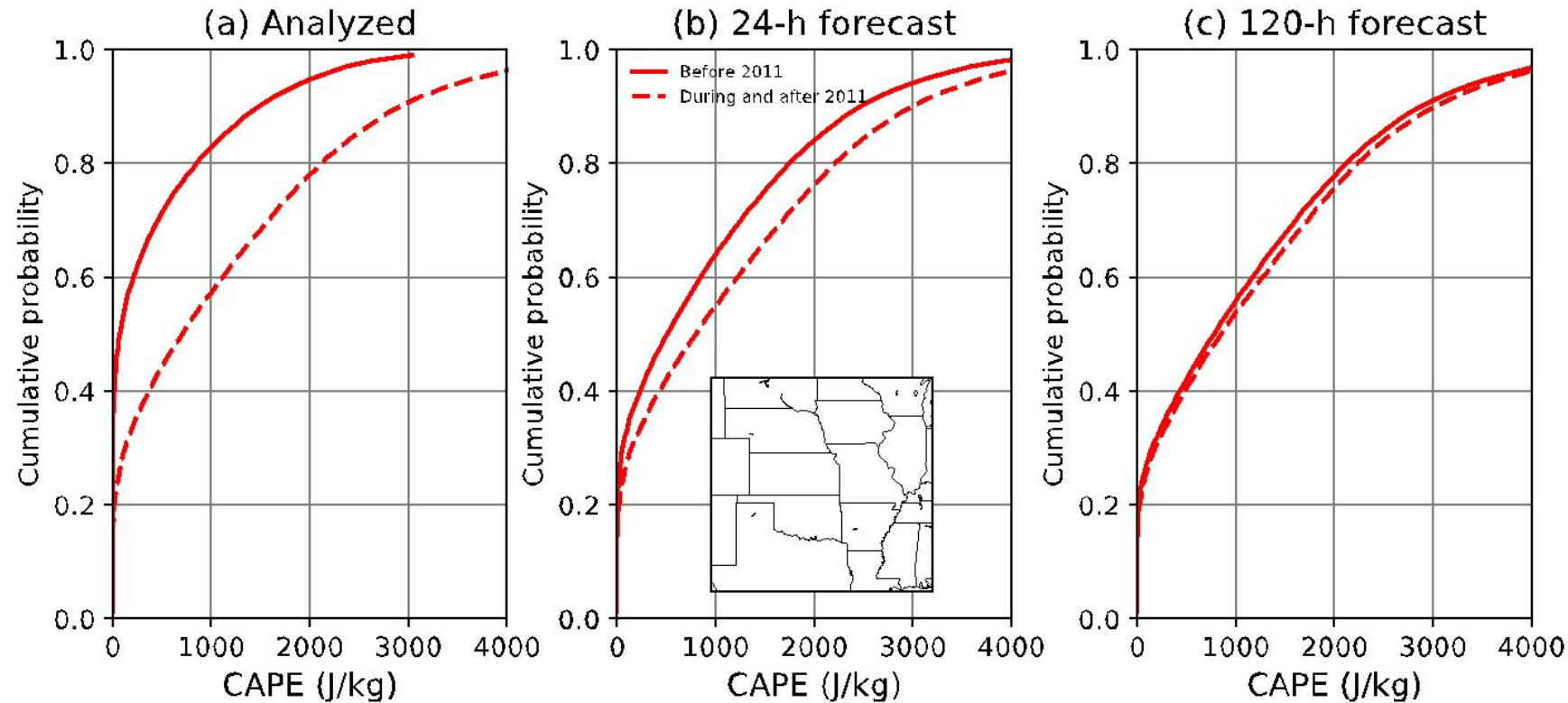
At least with simple methods  
(point-by-point logistic regression  
here), long time series of  
reforecasts are helpful.

Ref: Hamill et al. 2004 MWR, [here](#).



# Modern reanalyses are computationally expensive and difficult to produce. Why should we regularly perform them?

Apr-May-Jun CDFs of CAPE in GEFSv10 reforecasts 00 UTC



The previous-generation NOAA modern-era reanalysis, CFSR, changed from a fixed system prior to 2011 to the operational DA system after 2011. This illustrates the potential statistical inconsistency of DA on forecast.

# Optimizing for different users has led to different choices by reanalysis producers.

Users who want accurate  
analysis state estimates



“Climate”  
reanalysis  
system

- Attention to:
- (a) surface analyses.
  - (b) bias correction.
  - (c) coupling, perhaps.

Users who want reforecasts to  
be statistically consistent with  
real-time forecasts.



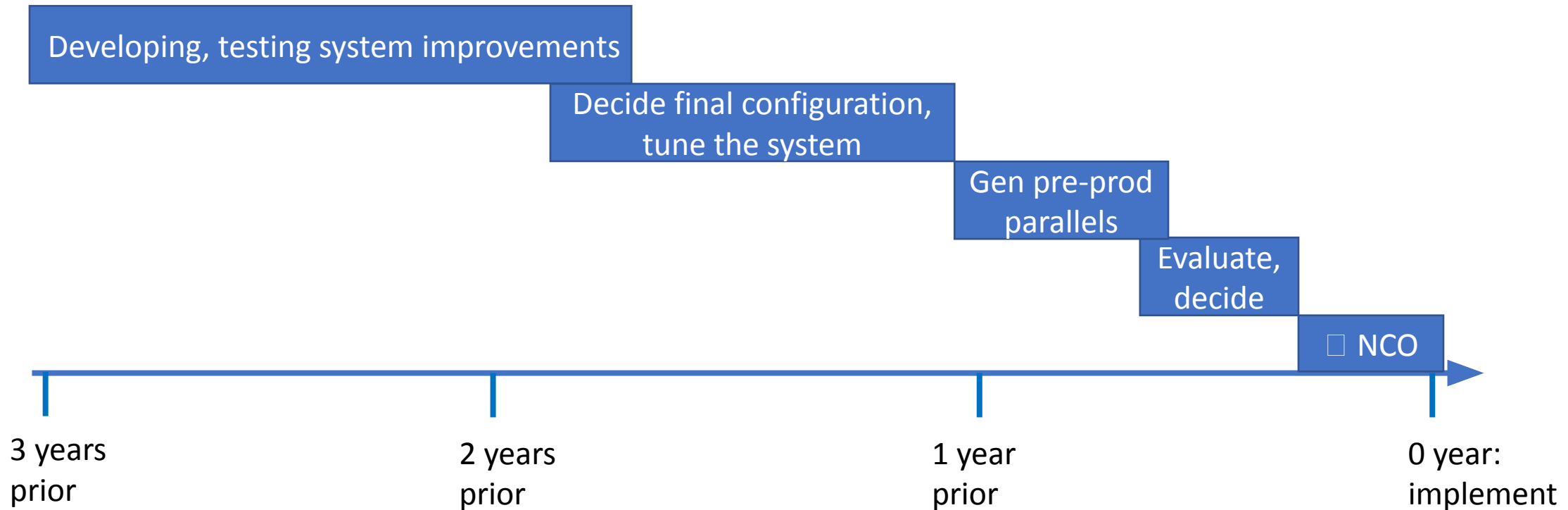
Reforecast  
initialization  
reanalysis  
system

Attention to consistency with  
operational DA methodology



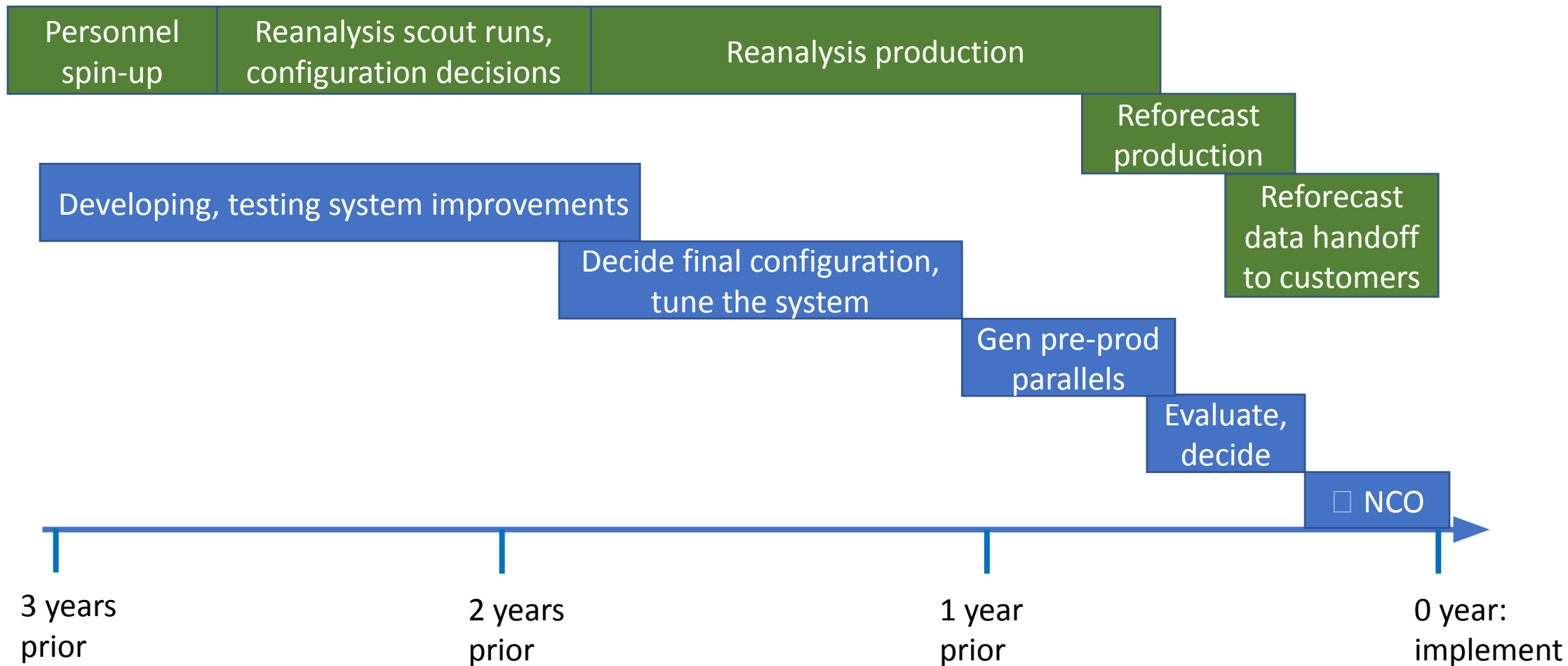
# A challenge with reanalysis production for reforecast init: synchronizing the reanalysis production with the operational upgrade schedule.

A rough guess at rolling out a next S2S prediction system in NOAA

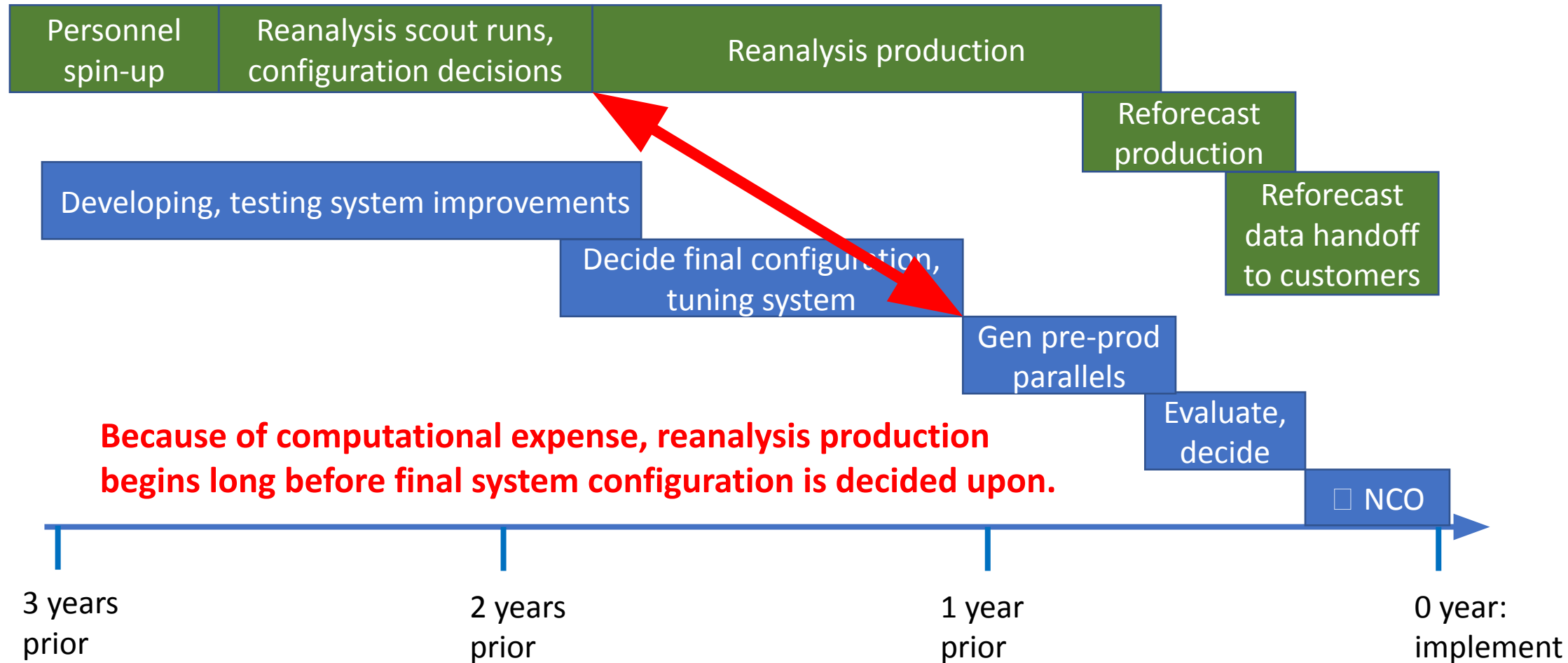


# A challenge with reanalysis production for reforecast init: synchronizing the reanalysis production with the operational upgrade schedule.

Because reanalysis production isn't institutionalized, like Copernicus



# A challenge with reanalysis production for reforecast init: synchronizing the reanalysis production with the operational upgrade schedule.



# More on reanalysis computational expense

- **NOAA reanalyses circa 2005:** uncoupled or weakly coupled, 3D-Var. Now comparatively cheap.
- **NOAA reanalyses circa 2020:** 80-member atmospheric En-Var. Computational expense and other factors necessitated:
  - DA uncoupled from land, ocean, sea ice.
  - Reduced-resolution EnKF provided background-error covariances.
  - Limited to 2000-2019 period.
- **Ideal future NOAA reanalysis.**
  - High spatial resolution, vertical and horizontal.
  - More observations assimilated, ideally consistent with other reanalysis producers to facilitate intercomparison.
  - Weakly or strongly coupled ocean, land, sea ice.
  - Consistent analysis increments applied to land from 2-meter observations.
  - Spanning decades prior to 2000, & advanced methods for dealing with data discontinuities.
  - DA improvements that facilitate use for climate analysis too (more later).
  - Include reprocessed observations, retrofitting new advances (e.g. like all-sky MW/IR) to older sensors

# Some possible ways of dealing with computational expense and long production period.

- **Surge reanalysis and reforecast production in the cloud** just prior to implementation.
  - But can we afford it? Will streams of 5 years cycled DA work in cloud?
- Find ways to **minimize computational expense**.
  - Keep resolution increases modest.
  - Computationally efficient machine-learning models to simulate forward models of the ocean, perhaps atmosphere.
  - Limit period of reanalysis.
- Use other reanalysis such as ERA5 adapted to be more statistically consistent with the operational DA – “**replay**” (like nudging). Then do full reanalysis more occasionally.



## Current generation instances

For the best performance, we recommend that you use the following instance types when you launch new instances. For more information, see [Amazon EC2 Instance Types](#).

Type	Sizes	Use case
DL1	dl1.24xlarge	Accelerated computing
F1	f1.2xlarge   f1.4xlarge   f1.16xlarge	Accelerated computing
G3	g3s.xlarge   g3.4xlarge   g3.8xlarge   g3.16xlarge	Accelerated computing
G4ad	g4ad.xlarge   g4ad.2xlarge   g4ad.4xlarge   g4ad.8xlarge   g4ad.16xlarge	Accelerated computing
G4dn	g4dn.xlarge   g4dn.2xlarge   g4dn.4xlarge   g4dn.8xlarge   g4dn.12xlarge   g4dn.16xlarge   g4dn.metal	Accelerated computing
G5	g5.xlarge   g5.2xlarge   g5.4xlarge   g5.8xlarge   g5.12xlarge   g5.16xlarge   g5.24xlarge   g5.48xlarge	Accelerated computing
G5g	g5g.xlarge   g5g.2xlarge   g5g.4xlarge   g5g.8xlarge   g5g.16xlarge	Accelerated computing

There are AWS instances that are likely suitable for running parallel reanalysis streams, GPU or CPU.



Type	Sizes	Use case
C4	c4.large   c4.xlarge   c4.2xlarge   c4.4xlarge   c4.8xlarge	Compute optimized
C5	c5.large   c5.xlarge   c5.2xlarge   c5.4xlarge   c5.9xlarge   c5.12xlarge   c5.18xlarge   c5.24xlarge   c5.metal	Compute optimized
C5a	c5a.large   c5a.xlarge   c5a.2xlarge   c5a.4xlarge   c5a.8xlarge   c5a.12xlarge   c5a.16xlarge   c5a.24xlarge	Compute optimized
C5ad	c5ad.large   c5ad.xlarge   c5ad.2xlarge   c5ad.4xlarge   c5ad.8xlarge   c5ad.12xlarge   c5ad.16xlarge   c5ad.24xlarge	Compute optimized
C5d	c5d.large   c5d.xlarge   c5d.2xlarge   c5d.4xlarge   c5d.9xlarge   c5d.12xlarge   c5d.18xlarge   c5d.24xlarge   c5d.metal	Compute optimized
C5n	c5n.large   c5n.xlarge   c5n.2xlarge   c5n.4xlarge   c5n.9xlarge   c5n.18xlarge   c5n.metal	Compute optimized
C6g	c6g.medium   c6g.large   c6g.xlarge   c6g.2xlarge   c6g.4xlarge   c6g.8xlarge   c6g.12xlarge   c6g.16xlarge   c6g.metal	Compute optimized
C6gd	c6gd.medium   c6gd.large   c6gd.xlarge   c6gd.2xlarge   c6gd.4xlarge   c6gd.8xlarge   c6gd.12xlarge   c6gd.16xlarge   c6gd.metal	Compute optimized
C6gn	c6gn.medium   c6gn.large   c6gn.xlarge   c6gn.2xlarge   c6gn.4xlarge   c6gn.8xlarge   c6gn.12xlarge   c6gn.16xlarge	Compute optimized
C6i	c6i.large   c6i.xlarge   c6i.2xlarge   c6i.4xlarge   c6i.8xlarge   c6i.12xlarge   c6i.16xlarge   c6i.24xlarge   c6i.32xlarge   c6i.metal	Compute optimized
D2	d2.xlarge   d2.2xlarge   d2.4xlarge   d2.8xlarge	Storage optimized



Product Details

<div>C6gC6gn</div>					
Instance Size	vCPU	Memory (GiB)	Instance Storage (GiB)	Network Bandwidth (Gbps)	EBS Bandwidth (Mbps)
c6g.medium	1	2	EBS-Only	Up to 10	Up to 4,750
c6g.large	2	4	EBS-Only	Up to 10	Up to 4,750
c6g.xlarge	4	8	EBS-Only	Up to 10	Up to 4,750
c6g.2xlarge	8	16	EBS-Only	Up to 10	Up to 4,750
c6g.4xlarge	16	32	EBS-Only	Up to 10	4750
c6g.8xlarge	32	64	EBS-Only	12	9000
c6g.12xlarge	48	96	EBS-Only	20	13500
c6g.16xlarge	64	128	EBS-Only	25	19000
c6g.metal	64	128	EBS-Only	25	19000
c6gd.medium	1	2	1 x 59 NVMe SSD	Up to 10	Up to 4,750
c6gd.large	2	4	1 x 118 NVMe SSD	Up to 10	Up to 4,750
c6gd.xlarge	4	8	1 x 237 NVMe SSD	Up to 10	Up to 4,750
c6gd.2xlarge	8	16	1 x 474 NVMe SSD	Up to 10	Up to 4,750
c6gd.4xlarge	16	32	1 x 950 NVMe SSD	Up to 10	4,750
c6gd.8xlarge	32	64	1 x 1900 NVMe SSD	12	9,000
c6gd.12xlarge	48	96	2 x 1425 NVMe SSD	20	13,500
c6gd.16xlarge	64	128	2 x 1900 NVMe SSD	25	19,000
c6gd.metal	64	128	2 x 1900 NVMe SSD	25	19,000

AWS C6g instances



# Instance purchasing options

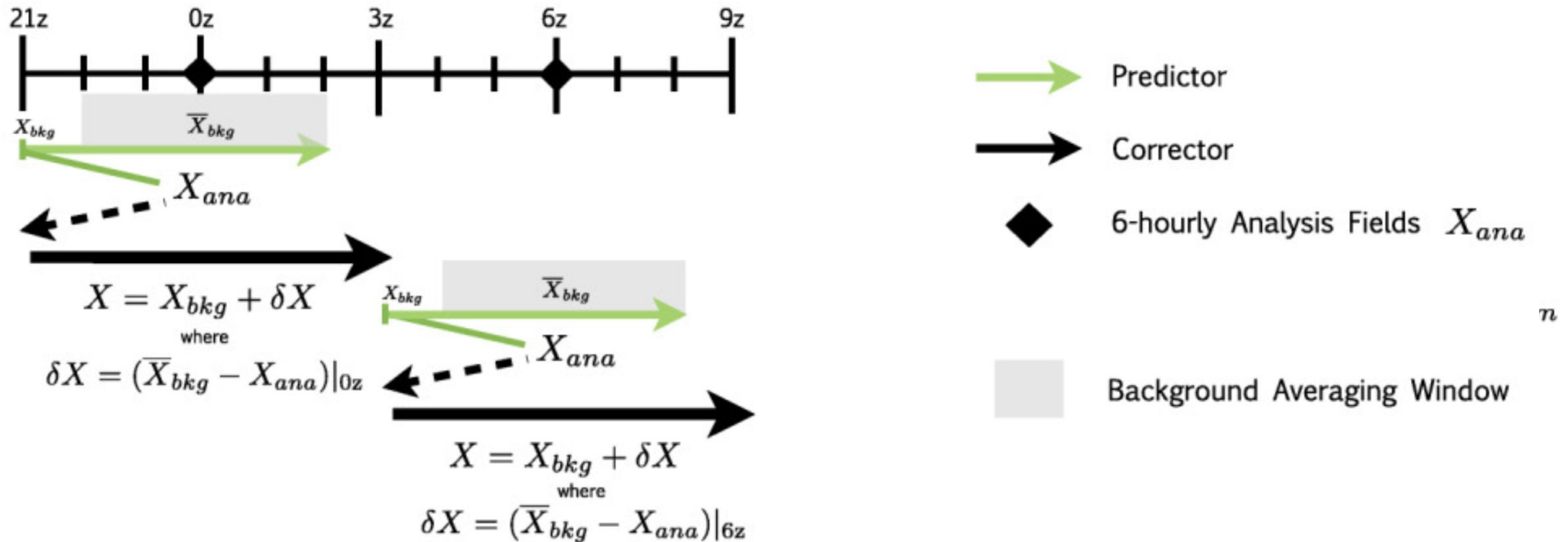
[PDF](#) | [Kindle](#) | [RSS](#)

Amazon EC2 provides the following purchasing options to enable you to optimize your costs based on your needs:

- **On-Demand Instances** – Pay, by the second, for the instances that you launch.
- **Savings Plans** – Reduce your Amazon EC2 costs by making a commitment to a consistent amount of usage, in USD per hour, for a term of 1 or 3 years.
- **Reserved Instances** – Reduce your Amazon EC2 costs by making a commitment to a consistent instance configuration, including instance type and Region, for a term of 1 or 3 years.
- **Spot Instances** – Request unused EC2 instances, which can reduce your Amazon EC2 costs significantly.
- **Dedicated Hosts** – Pay for a physical host that is fully dedicated to running your instances, and bring your existing per-socket, per-core, or per-VM software licenses to reduce costs.
- **Dedicated Instances** – Pay, by the hour, for instances that run on single-tenant hardware.
- **Capacity Reservations** – Reserve capacity for your EC2 instances in a specific Availability Zone for any duration.

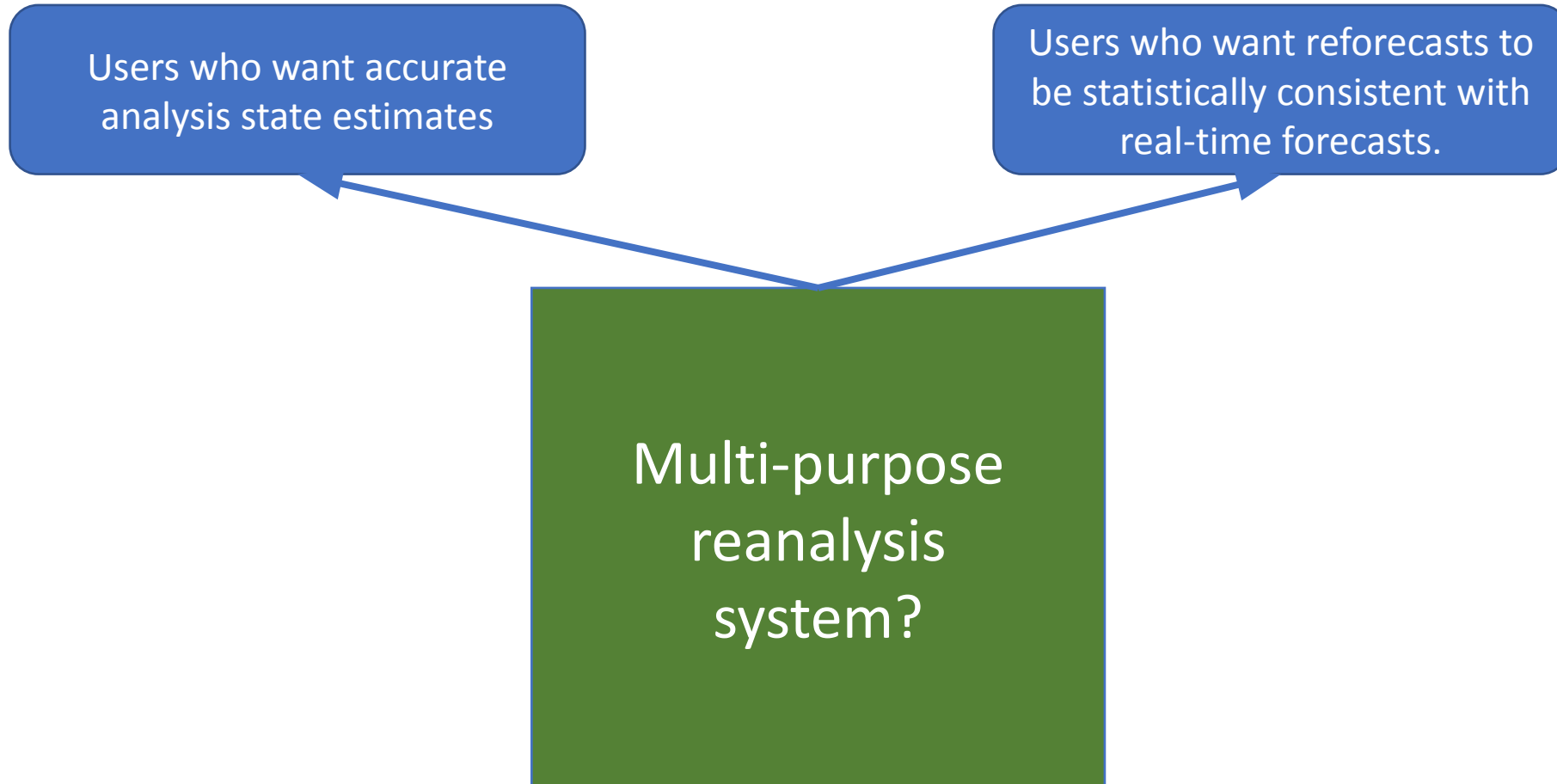
Spot instances are cheap, but would run sporadically. AWS has suggested that there are ways short of the spot market to minimize costs and guarantee throughput. Presumably similar for other cloud vendors.

# Replay (Orbe et al. 2017)



The replay technique uses the same Incremental Analysis Update (IAU) technique (Bloom et al., 1996) that was used to generate the MERRA-2 data assimilation, which consists of both predictor segments (green lines) and corrector segments (black lines). For (a) the RAnA simulation, a 5 h forecast centered about the analysis time (0z) is launched at 21z; an increment  $\delta X$  is then calculated as the difference between the time-averaged background state centered about 0z and a preexisting analysis field  $X_{ana}$  (black diamonds). The model is then backtracked to 21z and the increment  $\delta X$  is linearly applied to the background state over a 6 h corrector interval.

# One reanalysis product for multiple needs?



# Unification: incorporating characteristics of climate reanalysis into the operational DA.

- Weakly, quasi-strongly, or strongly coupled ocean, sea ice.
- Use of both atmospheric and surface states for  $H(\mathbf{x}^b)$  of surface-sensitive radiances for both ocean & atmosphere.
- Careful bias correction of background states.
- 2-m temperature and humidity analysis, 10-m wind analysis leveraging surface observations.
- Land assimilation: 2-m temperature and humidity increments make increments to land state.
- Coupled sea-ice assimilation.

# Conclusions

- NOAA is interested in modern-era reanalyses, but currently is optimizing for its primary application (reforecast initialization). This means not all customers are well-served with our reanalysis product.
- A long-term vision should be to build toward a reanalysis system that serves multiple purposes. Resource this!
- Reanalyses engineered for quality and completeness (high resolution, coupled, spanning many decades) are very computationally expensive, and we still are in an era of making severe tradeoffs to manage that expense.

# Biases in atmospheric reanalysis

- Contributed by background forecasts' systematic errors.
  - And correction of background bias via weak constraint 4D-Var, other methods challenging.
- Contributed by observations.
  - If other complimentary observations not available to anchor, analyses may inherit observation bias.
- Contributed by lack of observations.
  - Optimizing analyses to retain information from both tropospheric and surface observations difficult; some centers don't assimilate surface data.
- Contributed by lack of state coupling, which introduces transients with biases back into atmospheric background.