

REVIEW ARTICLE

Ensemble forecasting: A foray of dynamics into the realm of statistics

Jie Feng^{1,2,3}  | Zoltan Toth⁴ | Jing Zhang⁵ | Malaquias Peña⁶

¹Department of Atmospheric and Oceanic Sciences and Institute of Atmospheric Sciences, Fudan University, Shanghai, China

²Shanghai Key Laboratory of Ocean-land-atmosphere Boundary Dynamics and Climate Change, Shanghai, China

³Shanghai Academy of Artificial Intelligence for Science, Shanghai, China

⁴Global Systems Laboratory, NOAA, Boulder, Colorado, USA

⁵Shanghai Typhoon Institute, China Meteorological Administration, Shanghai, China

⁶Department of Civil and Environmental Engineering, University of Connecticut, Storrs, Connecticut, USA

Correspondence

Zoltan Toth, Global Systems Laboratory, NOAA, Boulder, Colorado, USA.
Email: zoltan.toth@noaa.gov

Funding information

National Natural Science Foundation of China, Grant/Award Numbers: 42105054, 42288101

Abstract

Uncertain quantities are often described through statistical samples. Can samples for numerical weather forecasts be generated dynamically? At a great expense, they can. With statistically constrained perturbations, a cloud of initial states is created and then integrated forward in time. By now, this technique has become ubiquitous in weather and climate research and operations. Ensembles are widely used, with demonstrated value. The atmosphere evolves in a multidimensional phase space. Does a cloud of ensemble solutions encompass the evolution of the real atmosphere? Theoretically, random perturbations in high-dimensional spaces have negligible projection in any direction, including the error in the best estimate, therefore consistently degrading that. As the bulk of the perturbation variance lies in the null space of error, samples in multidimensional space do not contain reality. An evaluation suggests that initial and short-range forecast error and ensemble perturbations are random draws from a high-dimensional domain we call the subspace of possible error. Error in any initial condition is partly a result of stochastic observational and assimilation noise, while perturbations explore other, mostly independent directions from the subspace of possible error that may have resulted from other configurations of stochastic noise. What benefits may arise from the deterministic projection of such noise? Consistent with theoretical expectations, ensemble members consistently degrade the skill of the unperturbed forecast until medium range. The mean and all other products derived from ensembles suffer an 18-hour loss in forecast information. Since information is a sufficient statistic, any rational user can benefit more from the unperturbed, than from an ensemble of weather forecasts. Furthermore, case-dependent variations in the distribution or spread of ensembles have no impact on commonly used metrics. Can alternative, statistical applications provide comparable, or even higher-quality probabilistic and other products, at the fraction of the cost of running an ensemble?

KEYWORDS

degrees of freedom, dynamical forecasting, ensemble forecasting, forecast uncertainty, statistical estimation of error variance

Jie Feng and Zoltan Toth contributed equally to this study.

1 | INTRODUCTION

Weather forecasting is one of the greatest success stories of natural sciences (Bauer *et al.*, 2015; Bennett and Richardson, 1923). Drawing on the theory of dynamics and thermodynamics, in an abstract setting, numerical models replicate the larger, resolved-scale dynamics of the atmosphere (Charney, 1949; Kalnay, 2003). In numerical weather prediction (NWP), observations of the atmosphere are collected and fused into an estimate of the initial state, called an analysis. Numerical forecasts initialized from such analyses then attempt to capture the temporal evolution of the atmosphere by exploiting deterministic relationships in nature. Useful forecast skill now extends to 10 days lead time and beyond (Bauer *et al.*, 2015; Zhang *et al.*, 2019) – a feat unimaginable just decades ago.

Despite continual reductions in initial error over the decades, error still amplifies in the forecasts. Eventually errors reach a level comparable with that in states randomly chosen from the climatic distribution, at which point forecasts become useless (Lorenz, 1982). By now it is well understood that the loss of forecast skill is intrinsic to a large class of aperiodic deterministic systems called chaotic dynamical systems (Li & Chou, 1997; Lorenz, 1963; Mu *et al.*, 2004; Thompson, 1957). As it is not due to methodological problems, this loss of skill is unavoidable (Lorenz, 1963). Weather is predictable – but only for a finite period.

Nature unfolds along a unique path in time and three-dimensional space. NWP forecasts attempt to predict this evolution in a similar form, as a unique sequence of events. Especially at longer lead times a single-value forecast in itself, however, can be rather deceptive. Such forecasts do not indicate how large their error may be, and which part of their variance will match reality. This is a major challenge for weather forecasters and users alike as for optimal decision-making the level, and possibly the nature of uncertainty must be known in advance (Leutbecher & Palmer, 2008).

After a brief review of statistical alternatives (Section 2), we introduce the concept of ensemble forecasting, a dynamical approach to assessing forecast uncertainty, along with its current status and presumed benefits (Section 3). Specific methodologies considered in this study, such as forecast system attributes, some metrics of forecast performance, including an analysis of perturbations in multidimensional space, and the sources of forecast error are discussed in Section 4. Long-held assumptions about ensembles are revisited in Section 5, while some theoretical explanation of the experimental results is offered in Section 6. The paper ends with some conclusions and a discussion (Sections 7 and 8, respectively).

2 | STATISTICAL METHODS

2.1 | Sampling

Statistical tools are available to describe uncertain quantities like weather analyses or forecasts. A sample or a distribution representing the expected error in the best estimate can readily show the range of values a quantity might take. Assuming, as an example, that the error in an analysis follows a normal distribution with known parameters, in Figure 1a the black curve centered around reality, whose exact value is unknown, indicates the possible position of an analysis, while the blue curve offers an example for a distributional estimate of reality, which, if the distribution is statistically reliable (i.e., perturbation variance equals error variance), is identical to the distribution of possible analyses, except translated to center on an arbitrarily selected realization of the analysis.

Error, by definition, is unknown at the time a forecast is made. Error variance, however, may be statistically assessed and used as an indicator of forecast uncertainty, as long as a representative joint forecast validation sample is available. Error variance in real time NWP guidance (i.e., analysis or forecast) fields at lead time i (\mathbf{G}_i), for example, can be estimated based on error variance in a sample of past guidance fields (\mathbf{S}_i) similar to \mathbf{G}_i in lead time, location, seasonality, regime, and so forth. (Hamill & Whitaker, 2006; Li & Ding, 2011; van den Dool, 1989; Zorita & von Storch, 1999):

$$e_{\mathbf{G}_i}^2 = E\left(e_{\mathbf{S}_i}^2\right), \quad (1)$$

where $E(\cdot)$ represents the expected value, and $e_{\mathbf{S}_i}^2$ and $e_{\mathbf{G}_i}^2$ are defined as:

$$e_{\mathbf{S}_i}^2 = |\mathbf{S}_i - \mathbf{T}|^2, \quad e_{\mathbf{G}_i}^2 = |\mathbf{G}_i - \mathbf{T}|^2, \quad (2)$$

and \mathbf{T} is the corresponding truth or its proxy (e.g., a verifying analysis, see Appendix A).

2.2 | Products

As an example, a set of surrogate or perturbed analyses or forecast fields (\mathbf{P}_i) can be created by the addition of perturbation fields ($\boldsymbol{\varepsilon}_i$) to the best, unperturbed single-value reference analysis or forecast field (sometimes also called “deterministic,” that from here on we call control, \mathbf{G}_i):

$$\mathbf{P}_i = \mathbf{G}_i + \boldsymbol{\varepsilon}_i, \quad (3)$$

Conveniently, past error patterns, if available, can serve as perturbations to create a sample of surrogate forecasts

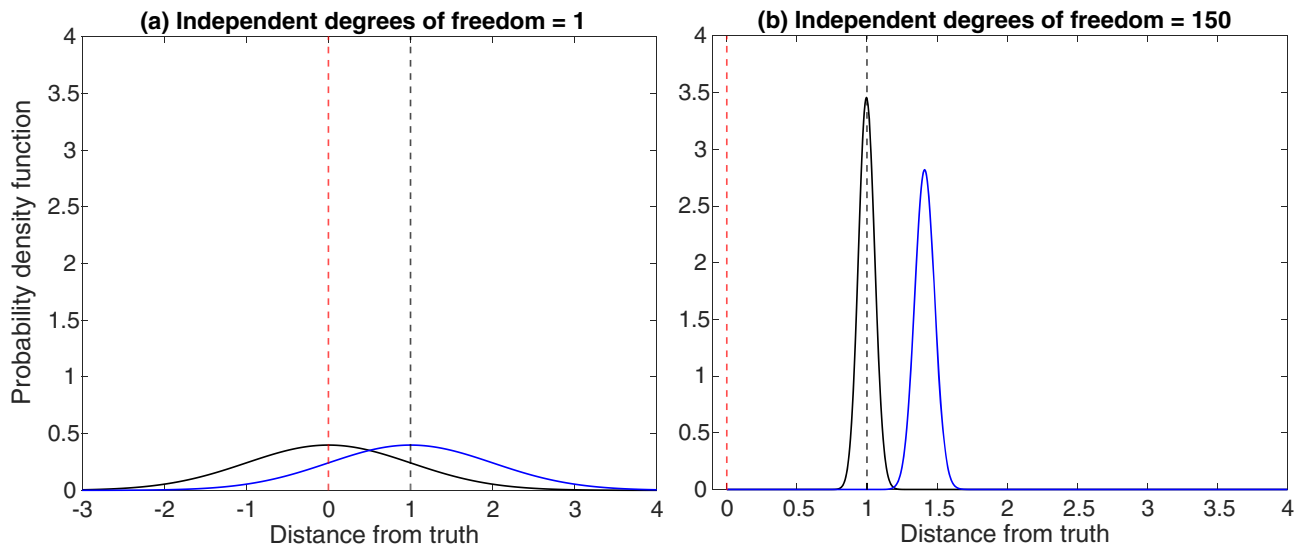


FIGURE 1 Reality (red dashed vertical line, unknown in practice), the distribution of its best estimate assuming error in it is known and normally distributed (black curve), an example for a best estimate (selected at a distance of the standard deviation from Reality, black dashed vertical line), and the estimated distribution of reality around the best estimate (blue curve) in a one- (panel a, directional distance) and 150-dimensional space (panel b, absolute distance). For further details, see text here and in Section 6.2.

(e.g., Delle Monache *et al.*, 2013). If a large enough archive of past error fields is not available, alternative sample generation methods include the addition of random noise (Leith, 1974; Palmer *et al.*, 1990), spatiotemporal shifts of a single forecast (neighborhood methods, e.g., Atger, 2001), or the collection of earlier initialized forecasts valid at the same time (lagged forecasts, Hoffman & Kalnay, 1983).

The mean of a sample is an often-used central tendency indicating the expected weather:

$$\mathbf{E}_i = \frac{1}{M_e} \sum_{k=1}^{M_e} \mathbf{P}_{i,k}, \quad (4)$$

where k is the index for perturbations, and M_e is the sample size. If perturbations are centralized before they are added to a reference state:

$$\sum_{k=1}^{M_e} \boldsymbol{\varepsilon}_{0,k} = 0, \quad (5)$$

the mean will equal the best estimate. In general, the mean captures the common component shared by all members. Typically, by filtering out presumably unpredictable noise, the mean of representative samples lowers forecast error.

To ensure statistical representativeness, the variance or spread of perturbation fields i is set equal to the estimated error variance in the best estimate:

$$V_i = \frac{1}{M_e} \sum_{k=1}^{M_e} |\mathbf{P}_{i,k} - \mathbf{E}_i|^2. \quad (6)$$

While the mean attempts to capture the predictable forecast signal, the spread (i.e., the standard deviation) measures the residual variability of sample points around their mean. Importantly, statistical sampling of forecast error involves the repeated, mechanistic insertion of perturbations around a single reference (control) forecast at every lead time (Figure 2a).

Using representative samples created around the best (control) single-value forecast, a variety of probabilistic and other products can be easily constructed in distributional or categorical (for semiclosed or closed intervals, see Anderson, 1996; Ebert, 2001) forms. For decades, statistical post-processing methods have been used to estimate and reduce forecast error, and generate well-calibrated forecasts in a variety of probabilistic and other formats (Scheuerer, 2014; Wilks, 2009). Due to limitations in methodology and the size of forecast archives, statistically generated surrogate forecasts, however, generally lack dynamical balance. Past forecast cases that best match the current forecast at a selected region and lead time, for example, lose such similarity at other locales and lead times. This is due to the high dimensionality of the atmospheric circulation (e.g., van den Dool, 1994). Hence to ensure representativeness, the selection of past forecast cases is often dependent on location and lead time (e.g., van den Dool, 1989), which results in perturbations that lack spatiotemporal and across-variable coherence or dynamical balance.

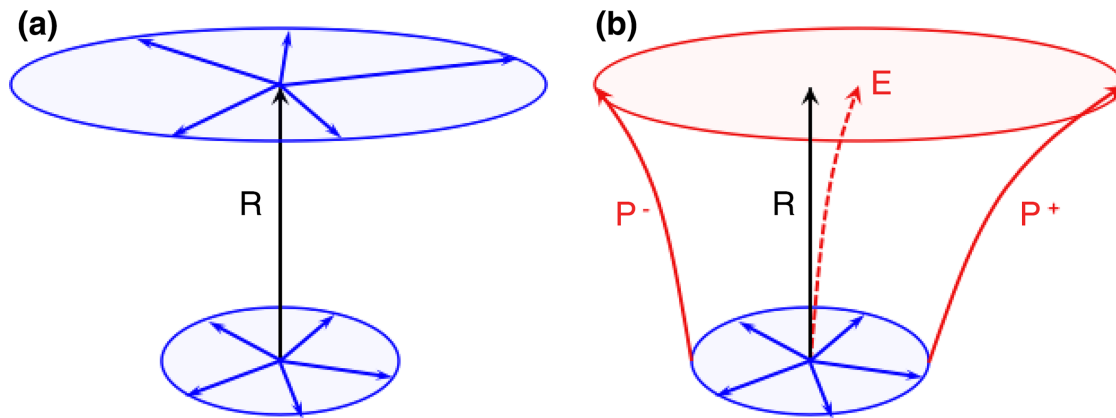


FIGURE 2 Schematic of statistical (a) vs dynamical (b) generation of forecast perturbations. In either case, initial perturbations (bottom ellipsoids) are centered on a reference initial condition (R, typically a control analysis and forecast, vertical black line). Forecast perturbations (top ellipsoids) are either statistically added and centered on R (a, blue arrows), or generated via the numerical integration of a dynamical model from perturbed initial conditions (b, red arrows). P^- , P^+ (red solid), and E (red dashed) represent two perturbations initially symmetric around, but later off-center of R, and the mean of the ensemble, respectively. For further explanation, see text.

3 | DYNAMICAL ALTERNATIVE

3.1 | Ensemble forecasting

Considering the limitations of statistical sampling algorithms and the success of the numerical approach to weather forecasting, a desire for the dynamical sampling of forecast uncertainty arose early on. Instead of the repetitive sampling of individual forecast variables (e.g., weather parameters at selected locations and lead times), why don't we sample the dynamical evolution of the entire atmosphere? In the 1960s an idea about a “glob of points, each of which would follow its own deterministic path” emerged (Edward Epstein, quoted by Lewis, 2005). The basic concept of ensemble forecasting is rather simple. Insert perturbations around the analysis of the atmosphere only once, at the initial time. To represent uncertainty in the analysis (Equation 1, see also Section 4.3.3), the magnitude of initial perturbations is set equal to that estimated in the analysis. And to retain skill in the mean, the initial sample is typically centered on the best, control estimate of the state (Equation 5). To create an ensemble, forecast perturbations are then dynamically generated by numerical integrations of the same (or to simulate model-related errors, a different, e.g., Houtekamer *et al.*, 2009) numerical model used to make the unperturbed control forecast (Figure 2b). A collection of such perturbed initial and forecast conditions are hence called an ensemble.

In the late 1980s and early 1990s, following experiments with models only about half the resolution of operational forecasts at the time, the idea gained momentum. In 1992, related efforts led to the operational implementation of the Global Ensemble Forecast System (GEFS)

at the National Centers for Environmental Prediction (NCEP, Toth & Kalnay, 1993). The routine weekend production of ensemble forecasts at the European Center for Medium Range Weather Forecasts (ECMWF) commenced shortly afterward (Molteni *et al.*, 1996). The rest is history (Lewis, 2005).

The dynamical generation of an ensemble, of course, comes at a significant cost. Depending on membership and resolution, in comparison with a single forecast, an order or two more computational resources may be required. Still, today dynamically generated ensembles constitute the main or sole mode of operation at most or all numerical weather and climate prediction centers (Chen & Li, 2020; Palmer, 2019; Zhou *et al.*, 2017). After decades of resistance, operational forecasters and other practitioners from a wide range of application areas (from hydrology, e.g., Schaake *et al.*, 2007, to agriculture, energy, and other sectors, e.g., Alemu *et al.*, 2011; Calanca *et al.*, 2011; Su *et al.*, 2014) and across many time-scales (from nowcasting, e.g., Liguori *et al.*, 2012, to multiseasonal and decadal forecasts, e.g., Krishnamurti *et al.*, 1999; Hou *et al.*, 2018; Liu *et al.*, 2023) have also embraced the practice (e.g., Bougeault, 2010). Ensembles and products derived from them, whether they represent the best possible guidance or not, are widely used, with proven value.

3.2 | Perceived benefits

Over the past decades, the potential benefits of ensemble forecasting have been discussed extensively. In this section we offer a brief overview of the perceived benefits. A more detailed analysis follows in Section 5.

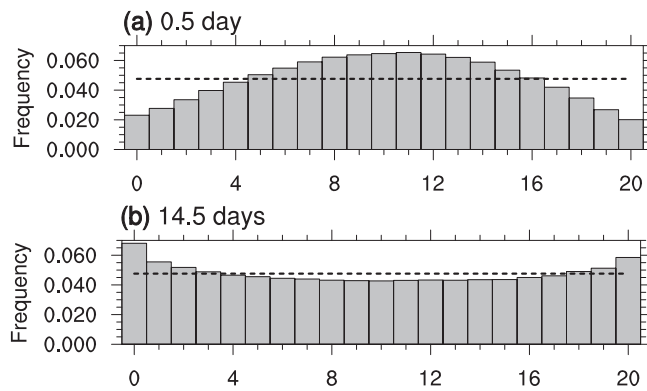


FIGURE 3 Talagrand (or analysis rank) diagram indicating the frequency of the verifying analysis falling into the intervals defined by the 20 ranked values of 500-hPa geopotential height ensemble member forecasts at individual grid points, aggregated over the northern hemisphere (NH) extratropics (30° – 65° N) over the three-month experimental period at 0.5 (a) and 14.5 days lead times (b). A flat distribution (dashed horizontal lines) indicates a perfectly reliable ensemble (where forecast probabilities of events exactly match their observed frequencies).

3.2.1 | Alternative scenarios

An attractive feature of ensembles is that they offer dynamically consistent alternative scenarios for future weather. Talagrand or Analysis Rank Histograms (Figure 3, C andille & Talagrand, 2005) demonstrate that the proxy for reality falls with about the same frequency in all intervals defined by an ordered set of ensemble members, an indication that ensemble scenarios are equally likely. A trivial but potentially powerful application is the direct feed of individual ensemble members into decision-making algorithms. A cost–benefit analysis in the context of the alternative forecast scenarios allows sophisticated users to optimize their weather-dependent course of actions (e.g., Alemu *et al.*, 2011; Khan *et al.*, 2021). A wide variety of probabilistic and other products can also be derived from such samples (Vannitsem *et al.*, 2021) just as easily as from statistical samples generated around single-value forecasts.

3.2.2 | Error reduction

Ensembles are well known for the low error in their mean (Equation 4). As shown in an example from the NCEP ensemble (Appendix A), the error in the mean (red line in Figure 4) is typically much below that in the control forecast run at the same resolution as the perturbed members (black). This is despite a noticeably higher error in the perturbed forecasts (blue line in Figure 4). The error reduction in the mean is considered a major benefit of ensembles. A

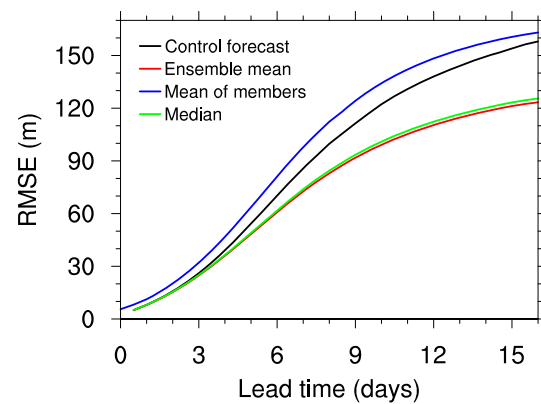


FIGURE 4 Perceived root-mean-squared error for the control (black), perturbed (blue), ensemble mean (red), and median (green) northern hemisphere (NH) extratropical 500-hPa height forecasts averaged over the three-month experimental period.

series of studies have suggested that the reduction in forecast error is dynamically conditioned, primarily due to a large projection of initial ensemble perturbations onto the “case-dependent” error in the control analysis (e.g., Toth & Kalnay, 1997, TK97, Ebert, 2001; Wei & Toth, 2003; Buizza *et al.*, 2008; Feng *et al.*, 2019). This presumed effect, often referred to as “nonlinear filtering,” is thought to “result in a superior ensemble mean forecast [compared] to a single or even higher-resolution control forecast” (Du, 2007). At the same time, it is maintained that a purely statistical “smoothing effect of [ensemble] averaging partially contributes to this superiority but ... in a much less degree ... compar[ed] to the nonlinear filtering” (Du, 2007).

3.2.3 | Spread–error relationship

Case-to-case variations in ensemble spread (Equation 6) are considered an important dynamical indicator of variations in expected forecast error variance (e.g., Buizza, 1997; Goerss, 2000; Murphy, 1988). Many link spatiotemporal variations in spread to fluctuations in atmospheric instabilities, presumably affecting forecast error variance (e.g., Ferranti *et al.*, 2015; Palmer, 2000). For further discussion, see Section 5.4.

3.2.4 | Probabilistic forecasts

A series of related papers (Christensen, 2015; Flowerdew, 2014; Hagedorn & Smith, 2009 and Roulston & Smith, 2003) compare verification scores for probabilistic products derived from an ensemble vs a higher-resolution control forecast. Roulston and Smith (2003), for

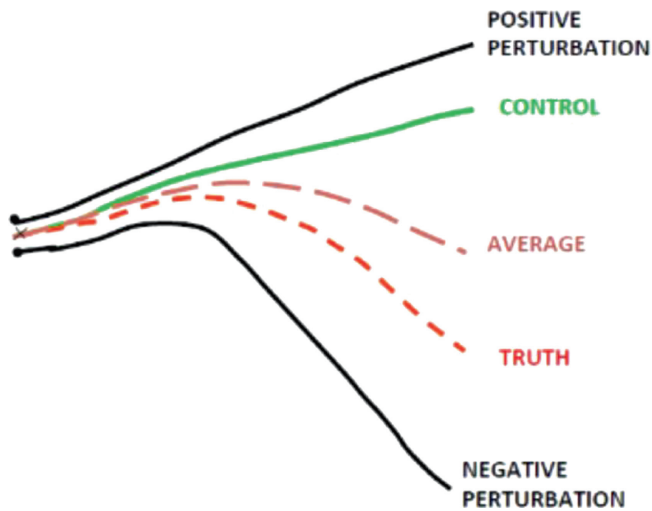


FIGURE 5 Schematic diagram of ensemble forecast trajectories: the control (green line), perturbed (black), and ensemble mean (long dashed brown) forecasts, and reality (dashed red). Courtesy of E. Kalnay, see text for details.

example, find that after applying very similar statistical post-processing methods, 3–10-day lead time ensemble-derived probabilistic forecasts have a much lower ranked probability score (RPS; Murphy, 1969) compared to products derived from a higher-resolution unperturbed forecast. Roulston and Smith (2003) and others attribute the favorable score for ensembles to their case-to-case varying distribution that provides “quantitative estimates of the likely forecast accuracy,” concluding that ensemble-based “prediction ... is inherently superior to a single “best guess” forecast.”

3.2.5 | Bracketing reality

From the beginning, a main objective of ensemble forecasting has been the dynamical sampling of uncertainty in the forecast evolution of the atmosphere. An ensemble brackets or encompasses truth if reality is contained in its range. As is well known, a statistically reliable M_e -member ensemble brackets any single indicator of reality or its proxy in the majority [i.e., $(M_e - 1)/(M_e + 1)$ fraction] of the cases (Descamps & Talagrand, 2007). As observed for commonly used variables, most of the time the proxy for truth (Appendix A) falls in the range of even somewhat unreliable ensemble forecasts (Figure 3a). Based on such experience in one dimension, the community has assumed that bracketing holds for the multidimensional space of atmospheric dynamics, too. This assumption is reflected in schematics like Figure 5 (reproduced from Kalnay, 2017), where the evolution of the real atmosphere is contained in, or dynamically bracketed by the cloud (i.e., the collection)

of ensemble forecast trajectories. This assumption will be evaluated in Section 5.5.

The introduction of ensembles was partly motivated by the applications, results, and expectations reviewed above. As ensembles proliferated in the weather forecast and user communities, some of the expectations solidified as presumptions. Many of these notions have never been critically examined. Motivated by, and building on the pioneering study of Leith (1974), the rest of this paper revisits some long-held assumptions about ensembles.

4 | CONCEPTS AND METHODOLOGY

The assessment of the quality of ensembles is critical to the optimal use and further development of forecast systems. In this section we review key concepts and tools we consider in the evaluation of ensemble forecasts.

4.1 | Forecast performance attributes

4.1.1 | Reliability

Based on a review of related literature, (Toth *et al.*, 2005) identified two forecast performance attributes: statistical reliability and statistical resolution (e.g., Murphy, 1972). Weather forecasts are in the form of abstract “signals,” each of which corresponds to a preferably unique weather event or condition in nature. Forecast symbols, as messengers in any communication, are arbitrary. Statistical reliability (e.g., Murphy, 1972) or calibration is one of two main attributes of forecast performance, assessing how truthful the forecast language is to its implied or expressly stated meaning. Specifically, reliability is not concerned about the sequence of forecast and observed events, just about their time average statistics. For example, is the mean of a sample of forecasts equivalent to the mean of corresponding observations? Naturally, metrics of reliability depend on the form of forecasts (i.e., symbols used, e.g., single-value or probabilistic, see Toth *et al.*, 2003). Therefore the reliability of forecasts expressed in different forms is quantitatively not comparable. Statistical reliability is key in the practical use of weather forecasts (Taillardat *et al.*, 2016). Fortunately, just as a text can be corrected for spelling errors without affecting its meaning, forecast bias can be statistically corrected based on past performance (i.e., calibration, e.g., Krzysztofowicz & Kelly, 2000).

4.1.2 | Resolution

Weather forecasts attempt to capture the temporal evolution of reality. As such, in contrast to their form, the sequence of events foreseen is the content of forecasts. Statistical resolution (e.g., Murphy, 1972) concerns how well the dynamical sequence of events in nature is captured by forecast signals indicating such events. In other words, resolution is a system's ability to foresee the sequence of future weather events, which in a loose sense can also be called the skill of forecast systems. Resolution is independent of the particular form or signals used and is arguably the inherent value, and the most critical attribute, of forecast systems. Note that resolution reflects only the similarity in the *sequence* but not in the long-term statistics of observed and forecast signals. As reliability is the other way around, the two main attributes of forecast systems are completely independent (Toth *et al.*, 2005).

Importantly, reliability and resolution are the only two attributes of forecast performance based on a comparison of forecasts and observations; other, diagnostic metrics concern only forecasts or observations alone. Therefore, our evaluation will focus on these two attributes since they provide a complete assessment of forecast system performance. Furthermore, as Krzysztofowicz (1992) noted, metrics of resolution are sufficient statistics in a sense that if their output can be calibrated, a forecast system with superior resolution will provide more economic benefit to any user compared to any other forecast system. So in terms of potential benefits, it is enough to compare the statistical resolution of forecast systems.

It follows that reliability and resolution of ensemble forecasts and products derived from them can (and preferably should) be evaluated separately. Most commonly used metrics of forecast performance, however, compound the two attributes with undetermined weights (and possibly include other elements, too). Since our focus is on forecast value, for a comparative evaluation of different forecast systems (such as single-value control, and multivalue ensemble forecasts), and for ease of interpretation, we will use a metric of resolution as a primary verification statistic.

4.2 | Forecast Information and Noise

As noted above, since they depend on the form of forecasts, reliability scores are quantitatively not comparable across different forecast systems. On the other hand, irrespective of their form, all forecast systems attempt to predict the sequence of future events; they differ only in what signals they use for communicating this. Unlike reliability, resolution therefore can be measured by common metrics, each assessing correlation between forecast and

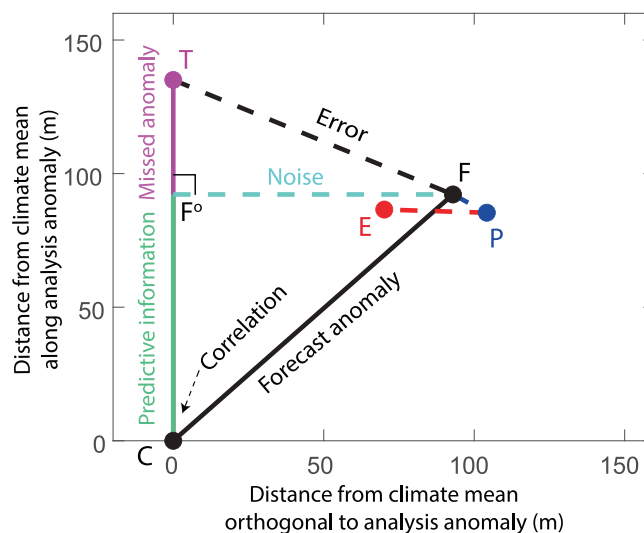


FIGURE 6 Schematic representing the phase space position of the seasonally and diurnally varying climatic mean (**C**), and unperturbed (**F**), perturbed (**P**), and ensemble mean forecast (**E**), and the corresponding truth or its proxy (**T**) on the Information (along the verifying analysis anomaly, vertical axis) and Noise (orthogonal to the verifying analysis anomaly, horizontal axis) plane. **P** and **E** are rotated into the **C**–**T**–**F** plane. Key performance metrics used in this study include the root-mean-squared error (**F**–**T**, or its square, the variance error, black dashed line); variances of forecast Information (**F**^o–**C**, solid green) and Noise (**F**–**F**^o, dashed cyan); the analysis anomaly missed by the forecast (**T**–**F**^o, solid pink); and pattern anomaly correlation (cosine of the angle at **C**). The position of the points indicates the performance of twice-daily eight-day lead time NCEP GEFS northern hemisphere (NH) extratropical (30°–65° N) 500-hPa height forecasts averaged over the December 1, 2017–February 28, 2018 experimental period. For further details, see text.

observed anomalies (Krzysztofowicz 1992, Krzysztofowicz & Evans, 2008).

4.2.1 | Information

In verification, we compare forecast quantities with the observed state, described here with its case-dependent anomaly from the climatic mean.ⁱ Let us consider forecast anomalies from the climatic mean of a model with realistic variability,ⁱⁱ standardized by the climatic variance (Figure 6). Further, we consider an orthogonal decomposition of forecast anomalies along, and orthogonal to the observed anomaly. Predictive capability or statistical resolution is measured here by the variance of the projection of forecast anomaly onto the observed anomaly, defined with respect to the climatic mean of nature:

$$I_i = \frac{|\mathbf{F}_i^o - \mathbf{C}|^2}{|\mathbf{T} - \mathbf{C}|^2}, \quad (7)$$

which we call forecast Information (I). \mathbf{F}_i and \mathbf{T} are an i lead time forecast and the corresponding truth or its proxy (e.g., a verifying analysis), respectively, \mathbf{C} is the climatic mean, \mathbf{F}_i^o is the orthogonal projection of \mathbf{F}_i on the observed anomaly $\mathbf{T} - \mathbf{C}$, and $|\bullet|$ is the Euclidean norm. In the rest of the manuscript, Information refers to I defined above. Information is the variance of the observed anomaly explained by a forecast, a direct measure of predictive capability. In other words, Information is the anomaly variance shared between reality and a forecast.

4.2.2 | Noise

In contrast, next we define variance in a forecast's anomaly that is orthogonal to the observed anomaly as Noise:

$$N_i = \frac{|\mathbf{F}_i - \mathbf{F}_i^o|^2}{|\mathbf{T} - \mathbf{C}|^2}. \quad (8)$$

Noise is an indicator for the level of divergence between a forecast and reality. Since Information that is identical to, and Noise that is unrelated to the observed anomaly constitute an orthogonal decomposition, for forecast systems with a realistic level of variance they are not independent quantities:

$$I_i + N_i = 1. \quad (9)$$

Information and Noise are therefore positively and negatively oriented, alternative and interchangeable metrics of forecast performance, respectively. Information/Noise variances standardized by the climatic variance range between 1/0 (perfect knowledge about nature) and 0/1 (no knowledge), respectively. Though related, Information and Noise defined above are different from "information entropy" (Shannon 1948) or noise used in signal processing (e.g., Tuzlukov, 2010, see Appendix B).

4.2.3 | Error

Error variance (Equation 1) is one of the most often-used metrics of forecast performance. Error measures the difference between a model forecast and reality. Theoretically, the initially quasi-exponential, then saturating growth of forecast error can be described by a logistic curve (Lorenz, 1982, see Appendix C). As seen from Figure 6 (dashed black line), error can be decomposed into Noise contained in (dashed cyan line, Equation 8), and Information missed by a forecast (continuous pink line, Equation C2).

4.2.4 | Information density

Pattern anomaly correlation (PAC or r_i , Jolliffe & Stephenson, 2003) is another commonly used performance metric, an inverse measure of the angle between forecast and verifying analysis anomalies taken from the climatic mean ($\mathbf{F}_i - \mathbf{C}$ and $\mathbf{T} - \mathbf{C}$, respectively in Figure 6). The square of PAC is interpreted here as Information density (I_i^d , see Figure 8d):

$$I_i^d = r_i^2 = \frac{I_i}{I_i + N_i}. \quad (10)$$

Note that for forecasts with the same, and only with the same anomaly variance, Information and Information density are interchangeable.

4.3 | Divergence of trajectory segments

At initial time, data assimilation systems capture partial Information about the state of nature, which numerical models then project into the future. Forecast error can be interpreted as the difference between segments of trajectories of dynamical systems. The difference between the evolution of two initially close segments on the trajectory of one, or two similar dynamical systems may be due to a number of factors.

4.3.1 | Difference in model dynamics

An important difference between the real atmosphere and its numerical models, beyond the latter being an abstract representation of reality, is that models explicitly consider only the larger-scale circulation. Following Leith (1974) and Zhou and Toth (2020), we assume that on larger scales well-resolved, numerical models replicate atmospheric dynamics in the extratropics near perfectly. Hence our study uses NH extratropical 500-hPa height as a primary dataset.

4.3.2 | Difference between equilibria

Though numerical models capture the dynamics of large-scale extratropical circulation well, their equilibrium (i.e., climatic mean) state differs from that of the real atmosphere. In other words, the attractor of numerical models is displaced from that of reality. When a model is initialized with a state close to that observed, it gradually drifts toward the model's own climatology. By definition, this process is governed by the stable dynamics of the model. As climatic drift in NH extratropical

500-hPa height is negligible, differences between the climatic mean of forecast and reanalysis fields are not considered in the definition of anomalies (Equations 7 and 8). Considering also Section 4.3.1, in the rest of this study we disregard the effect of model imperfections on forecast performance.

4.3.3 | Off-trajectory states

Though ideally the best (i.e., the control) and perturbed analyses of the atmosphere should all be in dynamical balance, in reality, both lie off the model trajectory (which approximates the trajectory of the large-scale motions of reality). Analysis fields contain random noise originating from both observational error and assimilation methods, while perturbations reflect intentionally imposed constraints (e.g., Tribbia & Baumhefner, 2004; Molteni *et al.*, 1996; TK97; Houtekamer & Mitchell, 1998). When a numerical model is applied to such imbalanced atmospheric states, over a relatively short (i.e., shorter than two-day) period, the stable part of dynamics pulls the evolving states close to the model trajectory. Once a forecast asymptotes the trajectory, the initial imbalance has no further effect on the divergence of trajectory segments. This is consistent with the findings that initial perturbations alter forecast performance just over a relatively short time period, after which only perturbation amplitude matters (Buizza *et al.*, 2005; Magnusson *et al.*, 2009; Raynaud & Bouttier, 2016). Therefore, imbalances in the evolution of error and perturbations are not considered explicitly in this study.

4.3.4 | Difference in the position on the trajectory

The divergence in the evolution of two points on the trajectory of a system that are originally close in phase space (but distant in time) is primarily driven by unstable dynamics (Buizza *et al.*, 1993; Feng *et al.*, 2018; Lorenz, 1982; Mu *et al.*, 2003). In the absence of model error, forecast uncertainty and the loss of predictability that ensembles aim to quantify arises due to such divergence. Assuming the variance distance between trajectory segments (i.e., error variance) follows a logistic evolution (see Section 4.2.3), both the growth of Noise and the loss of Information can be described analytically. As seen in Appendix C, due to the effect of unstable dynamics, with increasing lead time, Information is gradually converted into Noise variance, until all skill is lost. As the effects due to imbalances, or differences in the dynamics and equilibrium of systems are all negligible, error and perturbation behavior studied in

this paper are ascribed to the effect of unstable dynamics alone.

4.4 | Perfect model – perfect ensemble setup

Common verification practice also followed in this study involves the evaluation of forecasts such as the 20-member NCEP ensemble used in this study against verifying analysis fields. To eliminate the possible effect of specific data assimilation, modeling, and ensemble generation methods on evaluation results, in this study verification statistics will also be recalculated for 19 remaining members of the NCEP ensemble, replacing the verifying analysis with a randomly chosen ensemble member as truth. Reality and error in this simulated environment are generated by the same techniques as forecasts and perturbations, thus eliminating any influence from imperfect NWP methodologies. Following a long tradition established with the use of the term “perfect model” in observing system and other simulated experiments, we refer to this as a “perfect ensemble” setup. Note that the word “perfect” here does not imply an ultimate or ideal ensemble, but rather, a simulated environment where the ensemble forecast system uses a numerical model and perturbation generation method that are identical to those used in simulating reality and the error in its best estimate.

5 | EXPERIMENTAL RESULTS

Though ensembles are of multivalued form, just like a single control forecast, their members cover zero in probability space. Hence probabilistic and related products, just as from a single forecast, must be derived via statistical inter- and extrapolation (Vannitsem *et al.*, 2021). And whether single-value (e.g., Delle Monache *et al.*, 2013; Hamill & Whitaker, 2006; van den Dool, 1989) or ensemble-derived (e.g., Taillardat *et al.*, 2016), the reliability of probabilistic and other products can only be assessed and enforced by statistical methods, using a sample of past cases (Krzysztofowicz & Kelly, 2000).

Unfortunately, approximations in complex numerical models introduce biases into both single-value and ensemble forecasts. Difficulties in the estimation of the magnitude of initial, and in the representation of model-related errors also render the spread and distribution of ensemble forecasts unreliable (Vannitsem *et al.*, 2021). Hence in terms of one of the two major forecast performance attributes, statistical reliability, ensembles offer no benefit compared to single-value NWP forecasts. Products from both need to be statistically formulated, assessed and

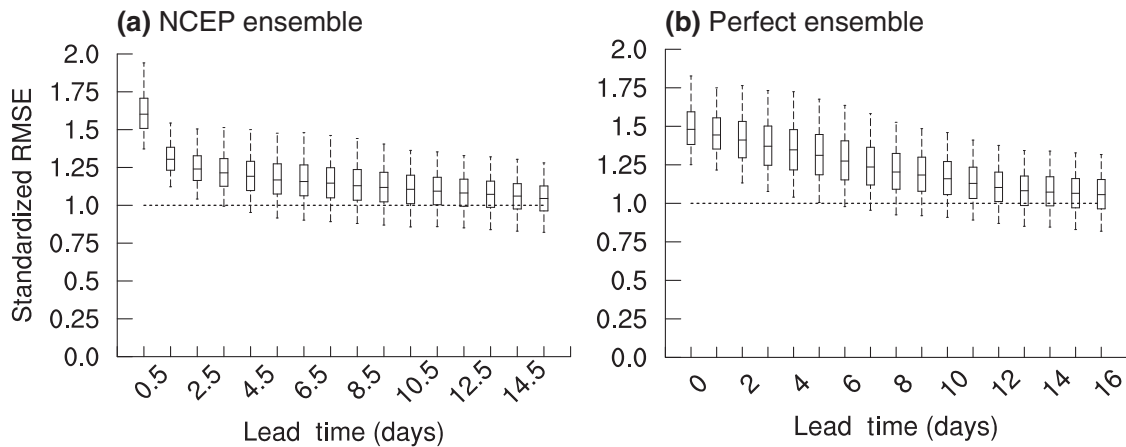


FIGURE 7 Northern hemisphere (NH) (30° – 65° N) 500-hPa height-perturbed forecast root-mean-squared error evaluated against the verifying analysis (a) and a randomly selected member (b), standardized by the error in 0.5–15.5 (panel a) and 0–16-day (panel b) control forecasts, ranked from lowest to highest, and averaged over all 180 cases. The top and bottom of whiskers and boxes represent the average of the extreme sample point and 25%/75% quantile values of the 20 and 19 ranked perturbed forecast error values in panels (a) and (b), respectively.

calibrated before their use. Next we evaluate what benefits ensembles may bring in terms of the second major forecast performance attribute, statistical resolution or forecast skill, or other unique aspects listed in Sections 3.2.1–3.2.5.

5.1 | Forecast quality

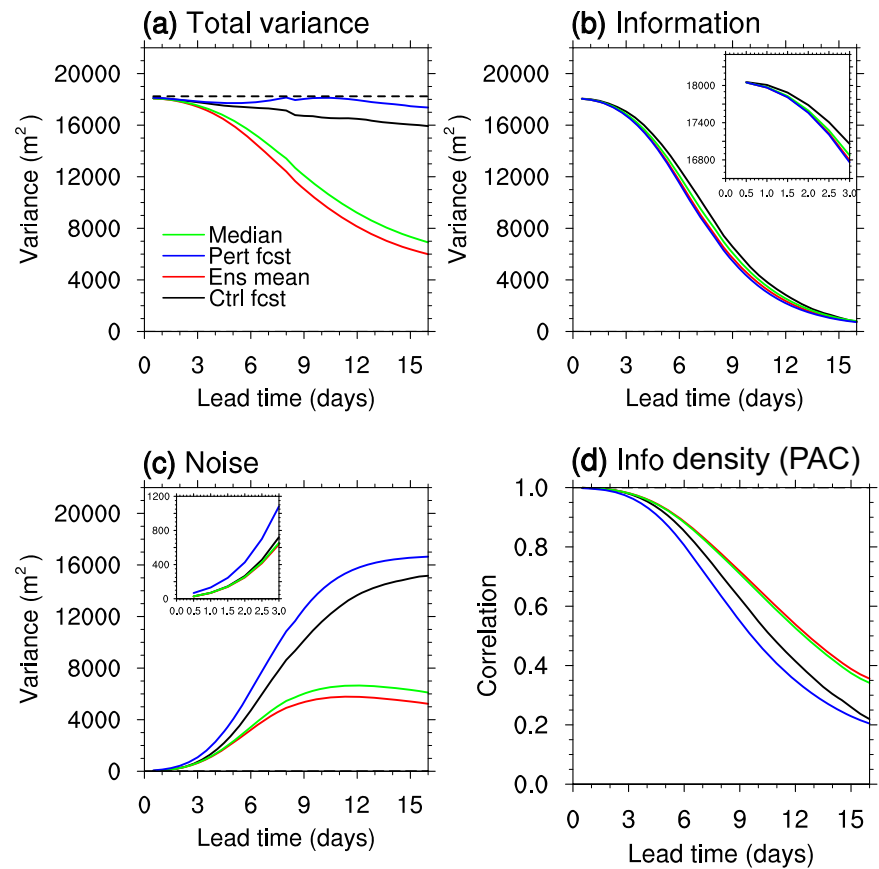
Diagrams like Figure 3b (Section 3.2.1) attest that members of ensembles offer equally likely scenarios. But are those scenarios also equally likely with the forecast started from the best, unperturbed control analysis?ⁱⁱⁱ An abundance of evidence indicates that they are not. Assuming perturbations are random draws from the distribution of initial error (Sections 2.1 and 3.1), based on simple statistical considerations the addition of perturbations to the best control analysis doubles their error variance compared to the control (Palmer *et al.*, 2006). This is born out in results from operational systems like the NCEP GEFS, where initial and short-range perturbed forecast error variance is about double that in the control forecast, negatively affecting performance at all ranges (cf. root-mean-squared [rms] error for the perturbed [blue] and control forecasts [black] in Figure 4). Moreover, we find that during the first few days, error variance in perturbed forecasts is higher not only in an expected sense, but also for each individual member. Shown in Figure 7a,b is the distribution of error in operational and perfect ensemble (see Section 4.4) perturbed forecasts, standardized separately in each case and for each lead time by the error in the control forecast. Apparently, shifts in phase space location introduced by ensembles induce a degradation in forecast quality similar to that due to spatiotemporal or other shifts made in a statistical sample generation.

Likewise, perturbed forecasts (blue line in Figure 8b) have lower Information compared to the control forecast (black line), reflecting an 18-hour loss in skill, equivalent to about an eight-year setback in NWP developments (Zhou & Toth, 2020). Significantly, the mean of the ensemble (red) shows a similar loss of Information.^{iv} The addition of random initial perturbations, like noise acquired in signal propagation, reduces Information in all members (not shown). One may argue, then, that unless other sources of Information are also considered, all products derived from ensemble forecasts will have Information lower than that in the control forecast, which is a key conclusion of this study. This is because new Information about nature cannot be created by taking a function of constituent members all characterized by lower quality. This situation is exemplified by the lower level of Information in the median of the ensemble (green curve in Figure 8b). We recall that Information is a measure of statistical resolution, or the inherent value in forecasts. Since Information is a sufficient statistic (Section 4.1), the results here indicate that any user may derive more benefit from a control forecast than from an ensemble. For optimal decision-making, one must use the control forecast, possibly with an added, statistically derived estimate of uncertainty. Is there some other value present in ensembles that may be missed by either of the two main forecast performance attributes, reliability or resolution?

5.2 | Error reduction

The lower error in the mean of an ensemble (cf. red and black lines in Figure 4) suggests yes, ensembles may have other benefits (Section 3.2.2). But how do we reconcile

FIGURE 8 Sample mean non-standardized (a) total variance, (b) Information variance, (c) Noise variance, and (d) Information density (or pattern anomaly correlation) of 500-hPa geopotential height forecasts in the northern hemisphere (NH) extratropics (30° – 65° N) over the 90-day experimental period (Appendix A). The dashed line in panel (a) indicates the climatic variance present in the analysis.



the reduction of error in the mean (Figure 4), a negatively oriented performance metric, with a concurrent decrease in Information (Figure 8b), a positively oriented metric? According to Equation (C3) (Appendix C), error can be reduced either by increasing Information, or decreasing Noise. As revealed by Figure 8c, the moderate reduction in Information is more than compensated with the large reduction of Noise in the mean. This is also apparent in the evaluation of eight-day forecasts in Figure 6. Apparently, the mean of an ensemble is a very efficient Noise filter. This is also reflected in the well-known, much smoother character of the mean as compared to single-value forecasts (Ansell, 2013), which is reflected in a significant reduction of overall variance in the mean (Figure 8a). Therefore, contrary to commonly held expectations (Section 3.2.2), error in the mean is reduced not because of a gain, but despite a loss of forecast Information, due to an effective reduction of unpredictable Noise.

5.3 | Probabilistic forecasts

Here we revisit the reason behind the lower error metrics found for ensemble- vs control-based probabilistic forecasts (Section 3.2.4). It turns out that unlike assumed by many, commonly used probabilistic scores like

Continuous Ranked Probability Score (CRPS, and its categorical equivalent, RPS) are not affected by variations in the shape or spread of forecast distributions (Hersbach, 2000). They depend only on the average of the spread of forecast distributions over the verification period. If not “case-dependent” variations in the shape of distributions, as suggested in the literature, then what explains the lower RPS error for ensemble-derived probabilistic forecasts? As Hersbach (2000) points out, CRPS (and hence RPS) is analogous to mean absolute error (MAE, which itself is closely related to error defined by Equation C1). The significantly lower RPS and other scores reported in Roulston and Smith (2003) and other studies for probabilistic forecasts derived from an ensemble vs a single control forecast is then a result of, just as in case of the error in the mean (Section 5.2), the reduced level of Noise in the position of ensemble distributions (i.e., their median) as compared to single-value forecasts (cf. green and black curves in Figure 8c).

5.4 | Spread–error relationship

An indication of the magnitude of forecast error (or error variance, Equation C1) by spatiotemporal fluctuations in ensemble standard deviation (or spread, Equation 6) is

another perceived benefit of ensembles (Section 3.2.3). The correlation between the two quantities, however, is rather low, explaining only about 10% in the day-to-day variability of the error magnitude (see, e.g., Figure 5 of Hopson, 2014). Perhaps not surprisingly, we found no anecdotal or documented evidence for the practical use of this relationship. As we saw in Section 5.3, fluctuations in spread certainly do not enhance forecast information. What may then explain the correlation between spread and error? By-and-large, the realizations of the atmosphere follow a multinormal distribution (Toth, 1995). In such a space, distances between states, just as in a univariate normal distribution, depend on a state's anomaly from the climatic mean (Li *et al.*, 2018). As forecast error measures the distance between trajectory segments of dynamical systems in multidimensional space (Section 4.3), it must also depend on the anomalies of the forecast and observed states. Evidence of this relationship for different forecast systems is presented by Toth (1991a, 1991b) and Kleeman (2011). We hypothesize that the weak relationship between spread and error may at least partially be explained by the dependence of both quantities on the climatic anomaly of the control forecast.

5.5 | Bracketing in multidimensional space

In one dimension, all perturbed states necessarily lie in the direction defined by reality and its forecast. The concept of bracketing, or encompassing truth is straightforward: reality must fall within the range of perturbed states (Section 3.2.5). Assuming a well-behaved unimodal distribution, this is possible only if some perturbed members have an error lower than the unperturbed estimate of reality, which is what we observe for all variables in today's ensembles. Does bracketing in any selected single direction guarantee bracketing in the multidimensional space of dynamics?

First we generalize the intuitive concept of bracketing into multidimensional spaces like that occupied by the dynamics of the atmosphere. There, just as in one dimension, bracketing is considered satisfied if reality falls in the range of perturbed states in the direction defined by reality and its forecast. This is the case-dependent direction of error in the unperturbed control, out of many independent degrees of freedom. Reduced error hence is still a necessary (but not sufficient) condition for bracketing in multiple dimensions. For bracketing to work in this space, perturbations must have a strong projection on, or be congruent with, the case-specific direction of error in the control. Bracketing case-specific error patterns in a multidimensional space is a much harder

challenge than bracketing single one-dimensional variables.

Experimental results in Figure 7 show that even for a subset of the atmosphere (500-hPa height variable over the NH extratropics) the necessary condition for bracketing of reduced error in the perturbed states is violated. Until day 3.5 and day 5, all members of the operational and perfect ensembles, respectively, have an error larger than that in the control forecast. An alternative interpretation of Figure 7 is that the time evolution of reality (or its proxy), the control forecast, and the range of perturbed forecasts are shown by the $Y=0$ line, the $Y=1$ line, and the boxplots, respectively. Figure 7 thus can be considered as a factual alternative to popular schematics like Figure 5 circulating in the community about ensemble forecasting. Clearly, in the case-dependent direction of error in the control forecast (and also in the control analysis in Figure 7b), reality or its proxy is far removed from the range of initial and short-range ensemble members. The widely-held assumption that the evolution of the atmosphere is contained in dynamically generated ensembles (Section 3.2.5) is untenable.

6 | THEORETICAL CONSIDERATIONS

6.1 | Simulation

In search of an explanation for the universal loss of skill, and the failure of dynamical bracketing demonstrated in Figure 7, we now turn our attention to the nature of high-dimensional spaces. For a quantitative assessment, we hypothesize that (i) unstable atmospheric dynamics responsible for the divergence of forecast and observed trajectory segments (cf. Section 4.3.4), as suggested by Toth (1991b, 1993) and Palmer *et al.* (2006), evolve in a multinormal space with a large number of independent and identically distributed (iid) variables^v (M_d), and that (ii) error and ensemble perturbations, after a short period of transitional behavior (see Section 4.3.3) are random draws from this domain we call the subspace of possible error. If these assumptions are valid, some basic features of forecast error and ensemble perturbation behavior should be statistically reproducible.

Our aim here is to compare the error in the initial unperturbed and perturbed states. While this can be accomplished for the perfect ensemble described in Section 4.4, we will use 12-hour forecasts instead as an indicator for error in the operational system. Plotted in Figure 9 are 12-hour lead time operational (20 dots, panel a) and perfect initial ensemble members (19 dots, panel b) for 180 cases along with the proxy for reality (vertical bar),

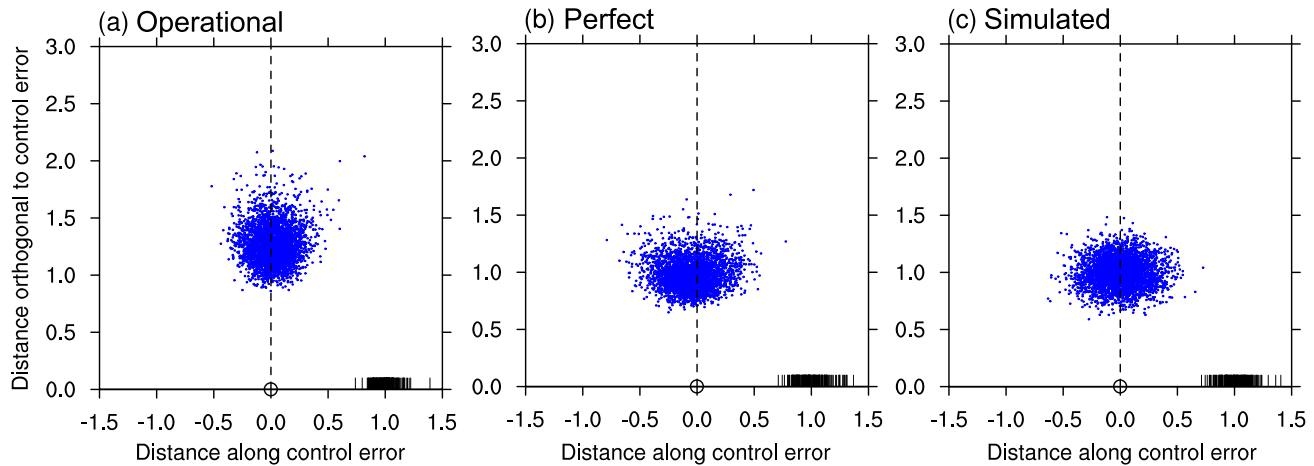


FIGURE 9 Perturbed northern hemisphere (NH) 500-hPa height (a) operational 12-hour forecasts, and (b) perfect initial ensemble members (3600 and 3420 individual blue dots from 20 and 19 members from each of 180 cases for panels a and b, respectively), plotted along (horizontal axis) and orthogonal (vertical axis) to the error in the unperturbed control 12-hour forecast (panel a, open circle at 0,0) or control analysis field (panel b) on a scale standardized by the sample mean error in the control, and the corresponding proxy for truth (black bars). Panel (c) is a statistical simulation of panel (b) with a 33 dof standardized multinormal distribution. For further details, see text.

as a function of distance from the control 12-hour forecast (panel a) or control analysis (panel b, both plotted at point 0,0) in the direction of error in the control (X axis, directional distance) and in the subspace orthogonal to it (Y axis, absolute distance in the null space of error in the control), on a scale standardized by the sample mean error in the control. Note that the distance of the bars and points from the control point (0,0) measures the size of error in, and perturbation around the control.

To validate the hypotheses above, we proceed with the generation of 20 random points from a distribution with a varying number of iid standardized normal variables (i.e., degrees of freedom [dof]). Just like in the perfect ensemble experiment reported in Figure 9b, one randomly chosen point is considered reality, while the remaining 19 the perturbed states. And just as is the case with the perfect ensemble, the simulation experiment is repeated 180 times. We find that the distribution of the error from the perfect and simulated ensembles are statistically indistinguishable at the 5% significance level for samples with a dof in the range of 28–38, with dof = 33 yielding the best fit (Appendix D), for which the results are plotted in Figure 9c.

Notable on all panels in Figure 9 is the small projection of perturbations introduced around the control forecast or analysis (0,0) onto the realization of error in the control (i.e., absolute value of X of perturbed points). This is in contrast with the magnitude of perturbations in the null space of error (i.e., Y value of perturbed points), which is comparable to the magnitude of error (i.e., distance of black

bars from the control at 0,0). Consequently, error variance for most members is almost doubled compared to the control (cf. the distance between reality and the control vs the perturbed states, consistent with rms error at 12-hour lead time in Figure 4). As error in all members is increased compared to the control, their cloud forms further away from reality. Consistent with Figure 7, reality or its proxy is not encompassed by either the operational or perfect ensembles. In all cases, the simulated ensemble also fails to bracket truth. Unlike in one dimension, statistical reliability (i.e., perturbation variance matching error variance) apparently does not imply bracketing in multidimensional space.

The remarkable visual and statistical similarity of the simulated (Figure 9c) to the perfect ensemble (Figure 9b), and a lesser, but still strong similarity to the operational ensemble data (Figure 9a) indicate that the experimental results are consistent with the hypotheses that (i) perturbation and error dynamics captured by NWP analyses and forecasts evolve in a multinormal space with a large number of iid variables, and that (ii) ensemble perturbations and error are indeed random samples from such a space. The space of resolved-scale error and perturbation dynamics is contingent on information captured in an analysis or forecast. The similarity of panels a and b in both Figures 7 and 9 also indicates that the problematic behavior observed in the operational ensemble, including their low skill and failure in bracketing cannot be addressed by perfecting data assimilation, modeling, or perturbation methodologies used.

6.2 | Interpretation

For an interpretation of error and perturbation results in Figure 9, we consider the orthogonal decomposition of anomaly variance into Information and Noise (Section 4.2). Depending on available observations and data assimilation techniques, NWP analyses and forecasts capture a certain amount of Information about the evolution of the larger-scale condition of the atmosphere, often considered deterministic. According to our hypotheses (Section 6.1), two states of the atmosphere given at the resolution of today's operational systems can differ in M_d independent ways, which for the NH extratropical height is estimated at 33. Given stochastic observational and methodological noise, error in any analysis is then just one random realization from this finer-scale "subspace of possible error." And perturbations which we assume are random draws from the same space simulate alternative realizations of analysis error that could have happened under different realizations of stochastic observational and methodological errors. Importantly, both Information about nature, and Noise contaminating a forecast are carried forward by the same model dynamics, albeit at different scales, used in numerical models.

In one dimension, reality and its best and perturbed estimates all occupy a single, common direction. Statistically reliable perturbations along this single direction bracket reality (Figure 1a). In one-dimensional space, statistical reliability is analogous with bracketing (Figure 3a). The dynamical evolution of the atmosphere, on the other hand, manifests in high-dimensional space. As the independent dof (M_d) increases, random draws from such a space spread out across more directions, lowering their expected projection on any single direction to $1/M_d$, including that of the error in unperturbed (control) estimates. Such behavior is often referred to as the "curse of dimensionality" (e.g., Bellman, 1961), which appears to be the fundamental cause of the failure of any sample, whether statistically or dynamically generated, in matching the level of Information in unperturbed estimates, or encompassing reality. As the bulk of perturbation variance projects into the null space of control error, error in each perturbed member is necessarily increased, failing to meet a necessary condition for bracketing.

A contributing factor to the loss of Information in, and the lack of bracketing by the perturbed members is the reduction of variability in the magnitude of both error and perturbations, which results in an even sharper separation of reality and its samples. Fluctuations in the magnitude of error (about 0.12 standardized units along X of the black bars in Figure 9b,c) and perturbations (0.12 along Y and 0.18 along X of the blue points) are greatly reduced compared to the standard deviation of 1 in one dimension. This

behavior is due to the Law of Large Numbers (e.g., Rose & Smith, 2002). If error in the best estimate (or control analysis) of a state is assumed to follow an iid normal distribution then theoretically, the distance of such guesses from reality follows a chi-squared distribution (black curve in Figure 1b), and the distance of perturbed states around any analysis from reality a non-central chi-squared distribution (blue curve in Figure 1b). The higher the dof, the narrower both of these distributions become. The demonstration in Figure 1b is for $M_d=150$ dof. Unlike in one dimension (Figure 1a), all perturbed states in high-dimensional spaces are further displaced from reality. Figure 1b hence indicates that the failure of operational, perfect, and simulated ensemble members to match the skill of unperturbed estimates, or to bracket reality is due to the peculiar geometry of high-dimensional spaces.

Finally, we contrast the time evolution of perturbations that are aligned, or congruent with, vs orthogonal to the error in the control. Initial perturbations congruent with the control error uniformly reduce or increase error in the perturbed state over the entire domain. Resulting perturbed states lie on a line defined by reality and the control initial condition. If error in each initial condition is assumed to grow logarithmically, the relative differences between smaller and larger initial errors will be retained in the forecast phase. Consequently, trajectories started with initial perturbations congruent with the control error will remain dynamically congruent in the forecast phase. Such forecast trajectories necessarily lie on a 2D surface defined by the trajectories of reality and the control forecast, ever diverging from, and never crossing each other (as suggested in Figure 5).

Such orderly error behavior is never observed with real-life ensembles. On the opposite, error curves for perturbed members evaluated over any subdomain display an incongruent, crisscrossing nature. This is evident in Figure 10, where the member that is best/worst over the NH extratropics in the 12–24-hour range (solid blue/red lines), for example, performs the worst/best a few days later (60–120 hours lead time range), respectively (or at other locales, not shown). This behavior can be explained by the random nature of initial perturbations in a high-dimensional space. With negligible projection on the actual error in the control, such perturbations improve/degrade the control initial condition in a random fashion over different parts of the domain. Model dynamics transposes the random initial spatial variations in skill into the time domain. The random fluctuations in forecast skill seen in Figure 10 hence arise as the influence of improvements and degradations in initial condition from different parts of the domain reaches the verification area. Clearly, ensemble perturbations behave like random (albeit spatiotemporally correlated) noise. And

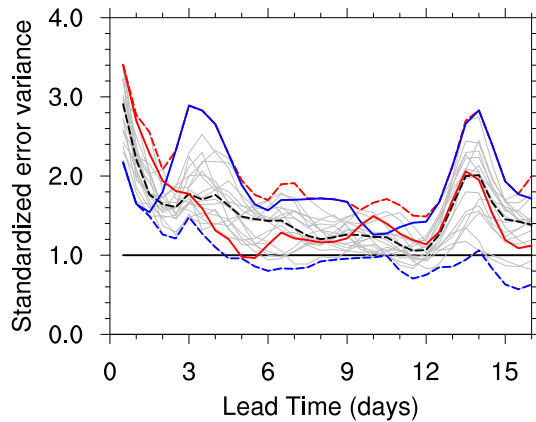


FIGURE 10 Same as Figure 7a, except error variance of individual forecasts against the verifying analysis for the single case initialized at 1200 UTC, 30 December 2017. The three dashed curves represent the error in the best (bottom, blue), median (middle, black), and worst member (top, red) at each lead time separately. The blue and red solid curves show the error variance in the members best and worst at the 12-hour lead time, respectively. Light gray curves show the error variance of individual members.

paradoxically, these random fluctuations provide statistical bracketing for single observed variables (Figure 3), while they fail to dynamically bracket the full state or its evolution (Figures 7 and 9).

6.3 | Nonlinear effects

To avoid the introduction of sampling error, initial perturbations are symmetrically arranged around the control analysis (Equation 5), setting the mean of operational ensembles equal to the control analysis. What explains the moderate and large reduction of Information and Noise in the ensemble mean compared to the control forecast, respectively? As noted by Gilmour *et al.* (2001), perturbations with amplitudes small relative to their saturation value develop quasi-linearly, leaving the mean mostly unaffected. This is reflected in the overlap of the black (control) and red (ensemble mean) curves in Figure 8b (Information), and especially in Figure 8c (Noise) in the 0–1-day lead time range.

Noticeable nonlinearities first emerge on the smallest scales due to the asymmetric evolution of the amplitude and position of affected features (Ancell, 2013). This results in a deviation of the mean from the control forecast. As nonlinear mixing in the 1–2-days lead time range is low, Noise removal is minimal; the difference between the mean and the control forecasts is dominated by the loss of Information. This is evidenced by the noticeably larger difference between the black (control) and red (ensemble

mean) curves for forecast Information (Figure 8b), compared to Noise (Figure 8c).

With increasing lead time, the phase and amplitude of perturbations on the smallest scales become fully randomized. At this stage of full nonlinear mixing, the mean of a typically sized ensemble removes a significant part of Noise present in the control forecast on scales with such fully saturated perturbation amplitudes. Simultaneously, error first in the control, then in the perturbed members also saturates, at which point all forecast Information on these small scales is lost. Due to the upscale propagation of energy, the same perturbation dynamics is repeated on successively larger scales. On scales with newly randomized perturbations, Information in the mean compared to the control forecast is temporarily reduced, after which a large part of Noise on such scales is removed. This succession of temporary reduction of Information and the additive removal of Noise on ever larger scales explains the increasing–steady–decreasing reduction of Information in the mean (compare black and red curves in Figure 8b), and the cumulative removal of Noise compared to the control forecast (compare black and red curves in Figure 8c), as a function of increasing lead time.

As perturbation energy moves upscale, the growth, as well as the overall variance of perturbations shifts to ever larger scales (cf. Figure 1 of Privé & Errico, 2015). This results in a general reduction of the independent degrees of freedom in perturbation dynamics, which explains the increase in the likelihood that the skill in some members over a limited domain rises above that of the control forecast, as noted earlier in Figure 7 for longer lead times. Ensembles, however, fail to bracket reality or its proxy at initial and short lead times even over a sub-domain of the 500-hPa height over the extratropical NH (with an estimated dof of 33). So bracketing observed at later lead times is only statistical, not dynamical in nature.

The dof of the full dynamics of short-range perturbations resolved by today's NWP systems is estimated to be in the range of 150–200 (see Appendix D). Bracketing in that space, as demonstrated for dof = 150 in Figure 1b, is even more challenging. Could the addition of more members help? The ratio of bracketing, more formally defined in Appendix E, is a function of dof and ensemble membership. In one dimension, the bracketing ratio with typical membership is sufficiently close to 1 $[(M_e - 1)/(M_e + 1)]$, ensuring that most of the time statistically reliable samples encompass a proxy of reality. In the high-dimensional space of the full resolved-scale dynamics of atmospheric circulation the chance of even large-size randomly generated ensembles encompassing reality, however, is astronomically low (see Figure E1).

7 | CONCLUSIONS

We exploit an orthogonal decomposition of forecast anomaly from the climatic mean into Information identical, and Noise orthogonal to the observed anomaly. Generally, Information about the state of natural systems is limited. Information is further reduced as forecast variance in chaotic systems like the atmosphere is gradually converted into Noise (Figure 6). For decades, statistical sampling has been successfully applied to assess uncertainty in weather forecasts (Figure 2a). Could forecast samples be generated dynamically, asked forerunners of ensemble forecasting. The practice of ensemble forecasting matured in the 1990s. Initial perturbations are added to the best estimate of the state, from which alternative scenarios are dynamically projected into the future (Figure 2b). After statistical calibration, probabilistic and other products derived from ensembles are widely used today, with demonstrated value.

Ensembles are assumed to (i) encompass the evolution of the real atmosphere (Figure 5); (ii) capture case-dependent variations in forecast error; and (iii) provide higher-quality single-value (ensemble mean, Figure 4) and (iv) probabilistic guidance. With a combination of theoretical and experimental approaches, these assumptions have been revisited. Using a statistical analysis, first we found that the divergence of segments of observed and/or forecast trajectory segments, and hence error and perturbation dynamics reside in a high-dimensional (150–200 independent degrees of freedom, Appendix D) domain we call the subspace of possible error. This subspace is contingent on the larger-scale condition of the deterministically evolving atmosphere, which one may associate with a “case.” Theoretically, sample points from high-dimensional spaces have negligible projection in any preselected direction, including the error in any initial state. Consequently, unlike in one dimension (Figure 1a), sample points in high dimensions consistently degrade the quality of the best estimate, and also miss to encompass reality (Figure 1b).

Information captured by an analysis is determined by the sophistication of the observing, data assimilation, and modeling systems. Experimental results suggest that error and perturbations are random draws from the high-dimensional subspace of possible error (Figure D1). Error in initial conditions results from specific realizations of stochastic noise in observations and data assimilation procedures, while perturbations represent alternative realizations of possible error that may have realized under different configurations of stochastic noise. As in real time Information and Noise are inseparable, numerical forecasts project their sum, the total initial variance, into the future. What value may the dynamical

generation of forecast samples (i.e., ensembles) via the deterministic projection of alternative Noise realizations bring?

An analysis of an operational and a perfect ensemble reveals that as theoretically expected, but contrary to assumption (i) above, initial perturbations and ensemble forecasts do not contain the state and evolution of the atmosphere (Figures 9 and 7, respectively). Also as expected, out to medium range, all members of the operational and perfect ensembles have larger error and less Information than that in the unperturbed control forecast. Unlike in one dimension (Figure 1a), ensemble members do not provide any scenario that is closer to reality than the control; ironically, they explore instead different ways that the control can be degraded (Figure 1b). And importantly, contrary to assumption (iii), the mean and arguably (iv) all probabilistic and other products derived from ensembles have less Information than that in the control forecast (Figure 8b).

Incidentally, an analysis by Hersbach (2000) shows that variations in the distribution of ensembles do not even have an effect on commonly used verification metrics. While other studies, as an alternative to assumption (ii), suggest that the low-level correlation found between case-to-case variations in spread and error may be explained by each being influenced by the amplitude of forecast anomalies. In any case, how could random draws from the subspace of possible error have any predictive information about the specific realization of error that is driven by stochastic observational and data assimilation processes? After all, it is only the subspace of possible error from which perturbations are also drawn, but not any specific, stochastically driven realization from it that is “case”-dependent (on the well-known large-scale conditions).

Our diagnosis indicates that the smaller error in the mean, a well-known benefit of ensembles, is due to an efficient filtering of Noise (Figure 8c) compared to individual forecasts. The smoother nature of the median of ensemble distributions is also what explains the lower scores found in probabilistic forecasts derived from an ensemble vs a control. Unfortunately, nonlinear filtering removes not only Noise, but some forecast Information as well (Figure 8b). Interestingly, Information is preserved in the mean only during the early, linear phase of the evolution of perturbations where their initial symmetry is still preserved and where ensembles are generally considered useless. Later, the loss of Information in the mean and other products amounts to an about 18-hour loss of lead time in warning about future weather events, or an eight-year setback in international NWP developments. The significance of this is that since Information is a sufficient verification statistic, any rationally acting user benefits more

from an unperturbed control than from an ensemble of forecasts.

Importantly, all behavior observed in operational ensembles is reproduced with a perfect ensemble. This confirms that their failure to meet expectations is not due to methodological shortcomings but lies rather in the multidimensional and nonlinear nature of atmospheric dynamics. At a great computational expense, ensembles recreate the same Information present in the control forecast M_e times, albeit at a lower level, while with painstaking accuracy generating M_e alternative realizations of dynamically balanced error of a somewhat larger magnitude. Ensembles lack statistical reliability or any discernible benefit from case-dependent variations, and have demonstrably less Information. Should the use of statistical alternatives be reconsidered? Filtering applications may reduce Noise in the best estimate while preventing or flexibly controlling the loss of forecast Information. With developing machine-learning applications like recent data-driven weather modeling (Bi *et al.*, 2023; Chen *et al.*, 2023), spatiotemporal and cross-variable covariances may also be induced into statistically generated perturbations. All the while calibrated probabilistic and other products of interest can be derived from statistical samples of error in past control forecasts, instead of dynamically generated ensembles.

8 | DISCUSSION

8.1 | Applicability

Though real-life results in this work are presented only with a single configuration, the NCEP GEFS, they arise out of general system characteristics. Specifically, (a) the phase space of all complex systems is high-dimensional, in which (b) nonlinear saturation randomizes perturbations on increasingly larger scales, reducing Information in the entire distribution. Therefore, the main conclusions of this study may in general be applicable to numerically created ensembles of high-dimensional multiscale dynamical systems. As an example, finer-scale processes resolved by increased resolution models are accompanied by higher degrees of freedom where, compared to synoptic scales, saturation of error happens faster, resulting in an even earlier onset of nonlinear perturbation behavior compared to what is found with the NCEP ensemble.

8.2 | Continuous approach

Whether forecast samples are represented in a quantized form of a finite sample (i.e., ensembles), or by a continuous

function (e.g., the Liouville Equations; Ehrendorfer, 2006), the underlying problems highlighted above remain the same. The loss of Information and the lack of bracketing therefore may equally affect continuous or quantized dynamical estimates of forecast uncertainty. In light of the availability of viable statistical alternatives, Leith's (1974) early assessment about ensembles may be applicable to continuous dynamical approaches as well: "sample sizes $M_e > 1$ will have to be justified on the basis of the detailed knowledge obtained . . ."

8.3 | Stochastic perturbations

Traditionally, the effect of finer-scale processes on motions explicitly resolved in numerical models is parameterized deterministically, conditioned on the resolved scales. More recently, with the intent of producing stochastic perturbations, random processes are inserted into some parameterization schemes. When such perturbations are added to forecast states during model integration, ensembles may become more reliable (e.g., Buizza *et al.*, 1999; Berner *et al.*, 2009), with reduced error in their mean (e.g., Sardeshmukh *et al.*, 2023). Just like initial perturbations, these random perturbations, however, also increase forecast Noise in individual members, and reduce Information both in the members and their mean. From a forecast Information perspective, one might consider stochastic perturbations as adding insult to injury sustained from the introduction of initial perturbations first.

8.4 | Data assimilation

Ensemble-derived products used in data assimilation describe covariances in the behavior of short-range forecast error. These products are based on ensemble forecasts issued at an earlier time and valid at the time of the analysis. As such, covariances have no forecast Information about the future state of the atmosphere; rather, they help find the best estimate of reality, given available observational data. Key limitations of ensembles such as the loss of Information or the lack of bracketing therefore do not affect data assimilation applications. Interestingly, the high dimensionality of the space of error discussed in Sections 5.5 and 6 has long been recognized in the context of ensemble-based data assimilation (e.g., "localization" algorithm of Szunyogh *et al.*, 2008).

8.5 | New elements

The degraded performance of all short-range perturbed forecasts compared to the control forecast evaluated over large areas is not a new finding (Palmer *et al.*, 2006).

Important implications such as the loss of Information in all derived products, and the failure of dynamically generated ensemble forecasts to encompass reality, however, have not been previously recognized. Neither have the lower error in the mean or in probabilistic forecasts derived from ensembles been attributed exclusively to Noise filtering, nor has the random nature of ensemble perturbations in the high-dimensional subspace of possible error, or the significance of the stochastic nature of error been recognized. Correspondingly, some basic characteristics of ensembles and the potential viability of alternative statistical sampling methods have for many remained elusive.

ACKNOWLEDGEMENTS

We sincerely thank Drs. Jian-Wen Bao and Jeff Beck of the Physical Sciences Laboratory and the Global Systems Laboratory (GSL) of NOAA, respectively, for their valuable comments on earlier versions of this report. The constructive comments of two anonymous reviewers greatly improved the presentation of the material. The second author (ZT) is grateful for many invigorating discussions he had over the years with members of the ensemble forecast and user communities. Discussions with Drs. Feifan Zhou (Institute of Atmospheric Physics), Duane Rosenberg (GSL), and the encouragement and continued support of Jennifer Mahoney and Kevin Kelleher, current and former Directors, and Curtis Alexander and DaNa Carlis, current and former Deputy Directors of GSL are also gratefully acknowledged. We also thank Dr. Jun Du of the Environmental Modeling Center, NCEP for his help with some of the references. The first author (JF) was partially supported by the National Natural Science Foundation of China (Grant Nos 42288101 and 42105054).

DATA AVAILABILITY STATEMENT

Numerical analysis and forecast data used in this study were provided by Dr Xiaqiong Zhou at the Environmental Modeling Center (EMC) of NCEP. The standard deviation and climate mean datasets used in some calculations were kindly provided by Drs Bo Yang and Suranjana Saha (EMC/NCEP). The data can be downloaded from NCEP's operational products inventory under GFS Ensemble Forecast System (GEFS) at <https://www.nco.ncep.noaa.gov/pmb/products/gens/>.

ENDNOTES

ⁱDifferences between points in phase space are independent of the choice of the reference point (or origin) used in defining possible coordinate systems. Here we adopt a convenient and often-used representation of atmospheric states through their anomalies from the climatic mean (e.g., Chen & Li, 2021).

ⁱⁱWe note that forecasts from operational prediction systems have a realistic level of variability, i.e., the overall variance in forecast ($F-C$) and verifying proxy for truth (i.e., analysis) anomaly fields ($T-C$) are near equal (see, e.g., less than 10% deviation between the solid and dashed black curves in Figure 8a, introduced later).

ⁱⁱⁱWhether a control analysis (and forecast) is produced by a data assimilation system in practice or not is immaterial. Whether the primary estimate of the state of nature is in single-value form around which an ensemble is introduced a posteriori, or an ensemble, the mean of which necessarily has a smaller error (Leith, 1974, first full paragraph in the left column of p. 411), a state with a superior estimate either exists, or can be identified, from which deviations of other estimates can be considered “perturbations.”

^{iv}A degradation in performance was first pointed out by Leith (1974) in the case of an ideal ensemble formed around reality, in comparison with a perfect control forecast.

^vWhether an orthogonal basis describing such a space can be determined in practice or not is irrelevant for our study; we are concerned only about the number of independent normal iid variates (dof) of this space.

ORCID

Jie Feng  <https://orcid.org/0000-0002-2480-2003>

REFERENCES

- Alemu, E.T., Palmer, R.N., Polebitski, A. & Meaker, B. (2011) Decision support system for optimizing reservoir operations using ensemble streamflow predictions. *Journal of Water Resources Planning and Management*, 137(1), 72–82.
- Ancell, B.C. (2013) Nonlinear characteristics of ensemble perturbation evolution and their application to forecasting high-impact events. *Weather and Forecasting*, 28(6), 1353–1365. Available from: <https://doi.org/10.1175/WAF-D-12-00090.1>
- Anderson, J.L. (1996) A method for producing and evaluating probabilistic forecasts from ensemble model integrations. *Journal of Climate*, 9(7), 1518–1530.
- Atger, F. (2001) Verification of intense precipitation forecasts from single models and ensemble prediction systems. *Nonlinear Processes in Geophysics, European Geosciences Union (EGU)*, 8(6), 401–417 hal-00331059.
- Bauer, P., Thorpe, A. & Brunet, G. (2015) The quiet revolution of numerical weather prediction. *Nature*, 525, 47–55. Available from: <https://doi.org/10.1038/nature14956>
- Bellman, R.E. (1961) *Adaptive control processes*. Princeton, NJ: Princeton University Press.
- Bennett, A.A. & Richardson, L.F. (1923) Weather prediction by numerical process. *The American Mathematical Monthly*, 30(1), 33. Available from: <https://doi.org/10.2307/2298921>
- Berner, J., Shutts, G.J., Leutbecher, M., & Palmer, T.N. (2009). A spectral stochastic kinetic energy backscatter scheme and its impact on flow-dependent predictability in the ECMWF ensemble prediction system. *Journal of the Atmospheric Sciences*, 66(3), 603–626. Available from: <https://doi.org/10.1175/2008jas2677.1>
- Bi, K., Xie, L., Zhang, H., Chen, X., Gu, X. & Tian, Q. (2023) Accurate medium-range global weather forecasting with 3D neural networks. *Nature*, 619, 533–538. Available from: <https://doi.org/10.1038/s41586-023-06185-3>

- Bougeault, P. (2010) The THORPEX interactive grand global ensemble. *Bulletin of the American Meteorological Society*, 91, 1059–1072. Available from: <https://doi.org/10.1175/2010BAMS2853.1>
- Buizza, R. (1997) Potential forecast skill of ensemble prediction and spread and skill distributions of the ECMWF ensemble prediction system. *Monthly Weather Review*, 125, 99–119.
- Buizza, R., Houtekamer, P.L., Toth, Z., Pellerin, G., Wei, M.Z. & Zhu, Y.J. (2005) A comparison of the ECMWF, MSC, and NCEP global ensemble prediction systems. *Monthly Weather Review*, 133, 1076–1097.
- Buizza, R., Leutbecher, M. & Isaksen, L. (2008) Potential use of an ensemble of analyses in the ECMWF ensemble prediction system. *Quarterly Journal of the Royal Meteorological Society*, 134, 2051–2066. Available from: <https://doi.org/10.1002/qj.346>
- Buizza, R., Miller, M. & Palmer, T.N. (1999) Stochastic representation of model uncertainties in the ECMWF ensemble prediction system. *Quarterly Journal of the Royal Meteorological Society*, 125, 2887–2908. Available from: <https://doi.org/10.1002/qj.49712556006>
- Buizza, R. & Palmer, T.N. (1998) Impact of ensemble size on ensemble prediction. *Monthly Weather Review*, 126(9), 2503–2518.
- Buizza, R., Tribbia, J., Molteni, F. & Palmer, T. (1993) Computation of optimal unstable structures for a numerical weather prediction model. *Tellus*, 45A, 388–407.
- Calanca, P., Bolius, D., Weigel, A.P. & Liniger, M.A. (2011) Application of long-range weather forecasts to agricultural decision problems in Europe. *The Journal of Agricultural Science*, 149, 15–22.
- Candille, G. & Talagrand, O. (2005) Evaluation of probabilistic prediction systems for a scalar variable. *Quarterly Journal of the Royal Meteorological Society*, 131(609), 2131–2150.
- Chakravarti, I.M., Laha, R.G., Roy, J. & Roy, J. (1967) *Handbook of methods of applied statistics*, Vol. I. Hoboken, NJ: John Wiley and Sons, pp. 392–394.
- Charney, J. (1949) On a physical basis for numerical prediction of large-scale motions in the atmosphere. *Journal of Meteorological Research*, 6, 371–385.
- Chen, J. & Li, X. (2020) The review of 10 years development of the GRAPES global/regional ensemble prediction. *Advances in Meteorological Science and Technology*, 10(2), 9–29. Available from: <https://doi.org/10.3969/j.issn.2095-1973.2020.02.003>
- Chen, L., Zhong, X., Zhang, F. & coauthors. (2023) FuXi: a cascade machine learning forecasting system for 15-day global weather forecast. *Npj Climate Atmospheric Science*, 6, 190. Available from: <https://doi.org/10.1038/s41612-023-00512-1>
- Chen, X. & Li, T. (2021) An improved method for defining short-term climate anomalies. *Journal of Meteorological Research*, 35(6), 1012–1022. Available from: <https://doi.org/10.1007/s13351-021-1139-2>
- Christensen, H.M. (2015) Decomposition of a new proper score for verification of ensemble forecasts. *Monthly Weather Review*, 143(5), 1517–1532. Available from: <https://doi.org/10.1175/MWR-D-14-00150.1>
- Delle Monache, L., Anthony Eckel, F., Rife, D.L., Nagarajan, B. & Searight, K. (2013) Probabilistic weather prediction with an analog ensemble. *Monthly Weather Review*, 141(10), 3498–3516.
- Descamps, L. & Talagrand, O. (2007) On some aspects of the definition of initial conditions for ensemble prediction. *Monthly Weather Review*, 135, 3260–3272.
- Du, J. (2007) How to evaluate the quality of an EPS and its forecasts? Part VI in uncertainty and ensemble forecast. p. 19–24, Science and Technology Infusion Lecture Series, NOAA's National Weather Service, Office of Science and Technology. Available online from: <https://www.nws.noaa.gov/ost/climate/STIP/uncertainty.htm>
- Ebert, E.E. (2001) Ability of a poor man's ensemble to predict the probability and distribution of precipitation. *Monthly Weather Review*, 129(10), 2461–2480. Available from: [https://doi.org/10.1175/1520-0493\(2001\)129<2461:AOAPMS>2.0.CO;2](https://doi.org/10.1175/1520-0493(2001)129<2461:AOAPMS>2.0.CO;2)
- Ehrendorfer, M. (2006) The Liouville equation in atmospheric predictability. In: Palmer, T. & Hagedorn, R. (Eds.) *Predictability of weather and climate*, Vol. 9780521848824. Cambridge, UK: Cambridge University Press, pp. 59–98.
- Feng, J., Ding, R.Q., Li, J.P. & Toth, Z. (2018) Comparison of nonlinear local Lyapunov vectors and bred vectors in estimating the spatial distribution of error growth. *Journal of the Atmospheric Sciences*, 75, 1073–1087.
- Feng, J., Li, J.P., Zhang, J., Liu, D.Q. & Ding, R.Q. (2019) The relationship between deterministic and ensemble mean forecast errors revealed by global and local attractor radii. *Advances in Atmospheric Sciences*, 36(3), 271–278.
- Feng, J., Toth, Z. & Peña, M. (2020) Partition of analysis and forecast error variance into growing and decaying components. *Quarterly Journal of the Royal Meteorological Society*, 146(728), 1302–1321.
- Ferranti, L., Corti, S. & Janousek, M. (2015) Flow-dependent verification of the ECMWF ensemble over the euro-Atlantic sector. *Quarterly Journal of the Royal Meteorological Society*, 141(688), 916–924. Available from: <https://doi.org/10.1002/qj.2411>
- Flowerdew, J. (2014) Calibrating ensemble reliability whilst preserving spatial structure. *Tellus, Series A: Dynamic Meteorology and Oceanography*, 66(1). Available from: <https://doi.org/10.3402/tellusa.v66.22662>
- Gilmour, I. & Smith, L.A. (1997) Enlightenment in shadows. In: Kadtko, J.B. & Bulsara, A. (Eds.) *Applied nonlinear dynamics and stochastic systems near the millennium*. New York: AIP, pp. 335–340.
- Gilmour, I., Smith, L.A. & Buizza, R. (2001) Linear regime duration: is 24 hours a long time in synoptic weather forecasting? *Journal of the Atmospheric Sciences*, 58, 3525–3539.
- Goerss, J.S. (2000) Tropical cyclone track forecasts using an ensemble of dynamical models. *Monthly Weather Review*, 128(4), 1187–1193. Available from: [https://doi.org/10.1175/1520-0493\(2000\)128<1187:TCTFUA>2.0.CO;2](https://doi.org/10.1175/1520-0493(2000)128<1187:TCTFUA>2.0.CO;2)
- Hagedorn, R. & Smith, L.A. (2009) Communicating the value of probabilistic forecasts with weather roulette. *Meteorological Applications*, 16(2), 143–155. Available from: <https://doi.org/10.1002/met.92>
- Hamill, T. & Whitaker, J.S. (2006) Probabilistic quantitative precipitation forecasts based on reforecast analogs: theory and application. *Monthly Weather Review*, 134, 3209–3229.
- Held, I. (2015) Small Earth, deep atmosphere, and hypohydrostatic models. Retrieved 4 Nov 2022 from GFDL Blog site: https://www.gfdl.noaa.gov/blog_held/65-small-earth-deep-atmosphere-and-hypohydrostatic-models/
- Hersbach, H. (2000) Decomposition of the continuous ranked probability score for ensemble prediction systems. *Weather and Forecasting*, 15(5), 559–570.
- Hoffman, R.N. & Kalnay, E. (1983) Lagged average forecasting, an alternative to Monte Carlo forecasting. *Tellus A*, 35A,

- 100–118. Available from: <https://doi.org/10.1111/j.1600-0870.1983.tb00189.x>
- Hopson, T.M. (2014) Assessing the ensemble spread–error relationship. *Monthly Weather Review*, 142(3), 1125–1142 Retrieved Mar 14, 2022, from <https://journals.ametsoc.org/view/journals/mwre/142/3/mwr-d-12-00111.1.xml>
- Hou, Z., Li, J.P., Ding, R.Q. & Feng, J. (2018) The application of nonlinear local Lyapunov vectors to the Zebiak–cane model and their performance in ensemble prediction. *Climate Dynamics*, 51, 283–304.
- Houtekamer, P., Mitchell, H. & Deng, X. (2009) Model error representation in an operational ensemble Kalman filter. *Monthly Weather Review*, 137, 2126–2143. Available from: <https://doi.org/10.1175/2008MWR2737.1>
- Houtekamer, P.L. & Mitchell, H.L. (1998) Data assimilation using an ensemble Kalman filter technique. *Monthly Weather Review*, 126, 796–811.
- Jolliffe, I.T. & Stephenson, D.B. (Eds.). (2003) *Forecast verification: a Practitioner's guide in atmospheric science*. Chichester: Wiley.
- Kalnay, E. (2003) *Atmospheric modeling, data assimilation and predictability*. Cambridge: Cambridge University Press.
- Kalnay, E. (2017) Historical perspective: earlier ensembles and forecasting forecast skill. Presentation at the ECMWF Annual Seminar 2017: Ensemble prediction: past, present and future. 11–14 Sep. 2017, Reading, England. https://www.ecmwf.int/sites/default/files/medialibrary/2017-03/AS2017_Programme.pdf
- Khan, A.N., Iqbal, N., Rizwan, A., Ahmad, R. & Kim, D.H. (2021) An ensemble energy consumption forecasting model based on spatial-temporal clustering analysis in residential buildings. *Energies*, 14(11). Available from: <https://doi.org/10.3390/en14113020>
- Kleeman, R. (2011) Information theory and dynamical system predictability. *Entropy*, 13(3), 612–649. Available from: <https://doi.org/10.3390/e13030612>
- Krishnamurti, T.N., Kishtawal, C.M., LaRow, T.E., Bachiochi, D.R., Zhang, Z., Williford, C.E. et al. (1999) Improved weather and seasonal climate forecasts from multimodel superensemble. *Science*, 285, 1548–1550.
- Krzysztofowicz, R. (1992) Bayesian correlation score: a utilitarian measure of forecast skill. *Monthly Weather Review*, 120(1), 208–219. Available from: [https://doi.org/10.1175/1520-0493\(1992\)120<0208:bcsaum>2.0.co;2](https://doi.org/10.1175/1520-0493(1992)120<0208:bcsaum>2.0.co;2)
- Krzysztofowicz, R. & Evans, W.B. (2008) Probabilistic forecasts from the national digital forecasts database. *Weather and Forecasting*, 23(2), 270–289. Available from: <https://doi.org/10.1175/2007WAF2007029.1>
- Krzysztofowicz, R. & Kelly, K.S. (2000) Bayesian improver of a distribution. *Stochastic Environmental Research and Risk Assessment*, 14(6), 449–470. Available from: <https://doi.org/10.1007/PL00009785>
- Leith, C.E. (1974) Theoretical skill of Monte Carlo forecasts. *Monthly Weather Review*, 102, 409–418.
- Leutbecher, M. & Palmer, T.N. (2008) Ensemble forecasting. *Journal of Computational Physics*, 227(7), 3515–3539. Available from: <https://doi.org/10.1016/j.jcp.2007.02.014>
- Lewis, M. (2005) Roots of ensemble forecasting. *Monthly Weather Review*, 133(7), 1865–1885.
- Li, J.P. & Chou, J. (1997) The existence of the atmosphere attractor. *Science China Earth Sciences*, 40, 215–224.
- Li, J.P. & Ding, R.Q. (2011) Temporal–spatial distribution of atmospheric predictability limit by local dynamical analogues. *Monthly Weather Review*, 139, 3265–3283.
- Li, J.P., Feng, J. & Ding, R.Q. (2018) Attractor radius and global attractor radius and their application to the quantification of predictability limits. *Climate Dynamics*, 51, 2359–2374. Available from: <https://doi.org/10.1007/s00382-017-4017-y>
- Liguori, S., Rico-Ramirez, M., Schellart, A. & Saul, A. (2012) Using probabilistic radar rainfall nowcasts and NWP forecasts for flow prediction in urban catchments. *Atmospheric Research*, 103, 80–95.
- Liu, T., Gao, Y., Song, X., Gao, C., Tao, L., Tang, Y. et al. (2023) A multi-model prediction system for ENSO. *Science China Earth Sciences*, 66, 1231–1240. Available from: <https://doi.org/10.1007/s11430-022-1094-0>
- Lorenz, E.N. (1963) Deterministic nonperiodic flow. *Journal of the Atmospheric Sciences*, 20, 130–141.
- Lorenz, E.N. (1982) Atmospheric predictability experiments with a large numerical model. *Tellus*, 34, 505–513.
- Magnusson, L., Nycander, J. & Källén, E. (2009) Flow-dependent versus flow-independent initial perturbations for ensemble prediction. *Tellus, Series A: Dynamic Meteorology and Oceanography*, 61(2), 194–209.
- Molteni, F., Buizza, R., Palmer, T.N. & Petroliagis, T. (1996) The ECMWF ensemble prediction system: methodology and validation. *Quarterly Journal of the Royal Meteorological Society*, 122, 73–119.
- Mu, M., Duan, W.S. & Chou, J.F. (2004) Recent advances in predictability studies in China (1999–2002). *Advances in Atmospheric Sciences*, 21, 437–443. Available from: <https://doi.org/10.1007/BF02915570>
- Mu, M., Duan, W.S. & Wang, B. (2003) Conditional nonlinear optimal perturbation and its applications. *Nonlin Processes Geophysics*, 10, 493–501.
- Murphy, A.H. (1969) On the “ranked probability score”. *Journal of Applied Meteorology*, 8, 988–989.
- Murphy, A.H. (1972) Scalar and vector partitions of the probability score: part I. two-state situation. *Journal of Applied Meteorology and Climatology*, 11, 273–282.
- Murphy, J.M. (1988) The impact of ensemble forecasts on predictability. *Quarterly Journal of the Royal Meteorological Society*, 114, 463–493.
- Palmer, T. (2019) The ECMWF ensemble prediction system: looking back (more than) 25 years and projecting forward 25 years. *Quarterly Journal of the Royal Meteorological Society*, 145, 12–24. Available from: <https://doi.org/10.1002/qj.3383>
- Palmer, T., Buizza, R., Hagedorn, R., Lawrence, A., Leutbecher, M. & Smith, L. (2006) Ensemble prediction: a pedagogical perspective. *ECMWF Newsletter*, 106(106), 10–17.
- Palmer, T.N. (2000) Predicting uncertainty in forecasts of weather and climate. *Reports on Progress in Physics*, 63(2), 71–116. Available from: <https://doi.org/10.1088/0034-4885/63/2/201>
- Palmer, T.N., Mureau, R. & Molteni, F. (1990) The Monte Carlo forecast. *Weather*, 45, 198–207.
- Peña, M. & Toth, Z. (2014) Estimation of analysis and forecast error variances. *Tellus A: Dynamic Meteorology and Oceanography*, 66, 1. Available from: <https://doi.org/10.3402/tellusa.v66.21767>
- Privé, N.C. & Errico, R.M. (2015) Spectral analysis of forecast error investigated with an observing system simulation experiment. *Tellus A*, 67, 25977.

- Raynaud, L. & Boultier, F. (2016) Comparison of initial perturbation methods for ensemble prediction at convective scale. *Quarterly Journal of the Royal Meteorological Society*, 142, 854–866. Available from: <https://doi.org/10.1002/qj.2686>
- Rose, C. & Smith, M.D. (2002) mathStatICA: Mathematical Statistics with Mathematica. In *Compstat*. Heidelberg, Germany: Physica-Verlag HD, pp. 437–442. Available from: https://doi.org/10.1007/978-3-642-57489-4_66
- Roulston, M.S. & Smith, L.A. (2003) Combining dynamical and statistical ensembles. *Tellus, Series A: Dynamic Meteorology and Oceanography*, 55(1), 16–30. Available from: <https://doi.org/10.1034/j.1600-0870.2003.201378.x>
- Sardeshmukh, P.D., Wang, J.A., Compo, G.P. & Penland, C. (2023) Improving atmospheric models by accounting for chaotic physics. *Journal of Climate*, 36, 5569–5585. Available from: <https://doi.org/10.1175/JCLI-D-22-0880.1>
- Schaake, J.C., Hamill, T.H., Buizza, R. & Clark, M. (2007) HEPEX—the hydrological ensemble prediction experiment. *Bulletin of the American Meteorological Society*, 88(10), 1541–1547. Available from: <https://doi.org/10.1175/BAMS-88-10-1541>
- Scheuerer, M. (2014) Probabilistic quantitative precipitation forecasting using ensemble model output statistics. *Quarterly Journal of the Royal Meteorological Society*, 140(680), 1086–1096. Available from: <https://doi.org/10.1002/qj.2183>
- Shannon, C.E. (1948) A mathematical theory of communication. *Bell System Technical Journal*, 27(3), 379–423. Available from: <https://doi.org/10.1002/j.1538-7305.1948.tb01338.x>
- Su, X., Yuan, H., Zhu, Y., Luo, Y. & Wang, Y. (2014) Evaluation of TIGGE ensemble predictions of northern hemisphere summer precipitation during 2008–2012. *Journal of Geophysical Research—Atmospheres*, 119, 7292–7310. Available from: <https://doi.org/10.1002/2014JD021733>
- Szunyogh, I., Kostelich, E.J., Gyarmati, G., Kalnay, E., Hunt, B.R., Ott, E. et al. (2008) A local ensemble transform Kalman filter data assimilation system for the NCEP global model. *Tellus, Series A: Dynamic Meteorology and Oceanography*, 60(1), 113–130. Available from: <https://doi.org/10.1111/j.1600-0870.2007.00274.x>
- Taillardat, M., Mestre, O., Zamo, M. & Naveau, P. (2016) Calibrated ensemble forecasts using quantile regression forests and ensemble model output statistics. *Monthly Weather Review*, 144(6), 2375–2393. Available from: <https://doi.org/10.1175/MWR-D-15-0260.1>
- Thompson, P.D. (1957) Uncertainty of initial state as a factor in the predictability of the large scale atmospheric pattern. *Tellus*, 9, 275–295.
- Toth, Z. (1991a) Estimation of atmospheric predictability by circulation analogs. *Monthly Weather Review*, 119(1), 65–72 Retrieved Mar 14, 2022, from https://journals.ametsoc.org/view/journals/mwre/119/1/1520-0493_1991_119_0065_eoapbc_2_0_co_2.xml
- Toth, Z. (1991b) Circulation patterns in phase space: a multinormal distribution? *Monthly Weather Review*, 119(7), 1501–1511 Retrieved Mar 14, 2022, from https://journals.ametsoc.org/view/journals/mwre/119/7/1520-0493_1991_119_1501_cpipsa_2_0_co_2.xml
- Toth, Z. (1993) Preferred and Unpreferred circulation types in the northern hemisphere wintertime phase space. *Journal of Atmospheric Sciences*, 50(17), 2868–2888 Retrieved Mar 21, 2022, from https://journals.ametsoc.org/view/journals/atsc/50/17/1520-0469_1993_050_2868_paukti_2_0_co_2.xml
- Toth, Z. (1995) Degrees of freedom in northern hemisphere circulation data. *Tellus*, 47A, 457–472.
- Toth, Z. & Kalnay, E. (1993) Ensemble forecasting at the NMC: the generation of perturbations. *Bulletin of the American Meteorological Society*, 74, 2317–2330.
- Toth, Z. & Kalnay, E. (1997) Ensemble forecasting at NCEP: the breeding method. *Monthly Weather Review*, 125, 3297–3318.
- Toth, Z., Talagrand, O. & Zhu, Y. (2005) The attributes of forecast systems: a framework for the evaluation and calibration of weather forecasts. In: Palmer, T.N. & Hagedorn, R. (Eds.) *Predictability of weather and climate*, Vol. 9780521848824. Cambridge University Press, pp. 584–595. Available from: <https://doi.org/10.1017/CBO9780511617652.023>
- Tribbia, J.J. & Baumhefner, D.P. (2004) Scale interactions and atmospheric predictability: an updated perspective. *Monthly Weather Review*, 132, 703–713. Available from: [https://doi.org/10.1175/1520-0493\(2004\)132<0703:SIAAPA>2.0.CO;2](https://doi.org/10.1175/1520-0493(2004)132<0703:SIAAPA>2.0.CO;2)
- Tuzlukov, V. (2010) *Signal processing noise, electrical engineering and applied signal processing series*. CRC Press, p. 688.
- van den Dool, H.M. (1989) A new look at weather forecast through analogs. *Monthly Weather Review*, 117, 2230–2247.
- van den Dool, H.M. (1994) Searching for analogues, how long must we wait? *Tellus A*, 46, 314–324. Available from: <https://doi.org/10.1034/j.1600-0870.1994.t01-2-00006.x>
- Vannitsem, S., Bremnes, J.B., Demaeyer, J., Evans, G.R., Flowerdew, J., Hemri, S. et al. (2021) Statistical postprocessing for weather forecasts review, challenges, and avenues in a big data world. *Bulletin of the American Meteorological Society. American Meteorological Society*, 102, E681–E699. Available from: <https://doi.org/10.1175/BAMS-D-19-0308.1>
- Wei, M. & Toth, Z. (2003) A new measure of ensemble performance: perturbations versus error correlation analysis (PECA). *Monthly Weather Review*, 131, 1549–1565.
- Wilks, D.S. (2009) Extending logistic regression to provide full-probability-distribution MOS forecasts. *Meteorological Applications*, 16(3), 361–368. Available from: <https://doi.org/10.1002/met.134>
- Zhang, F., Sun, Y.Q., Magnusson, L., Buizza, R., Lin, S.-J., Chen, J.-H. et al. (2019) What is the predictability limit of midlatitude weather? *Journal of the Atmospheric Sciences*, 76, 1077–1091. Available from: <https://doi.org/10.1175/JAS-D-18-0269.1>
- Zhou, F. & Toth, Z. (2020) On the prospects for improved tropical cyclone track forecasts. *Bulletin of the American Meteorological Society*, 1–55, E2058–E2077. Available from: <https://doi.org/10.1175/BAMS-D-19-0166.1>
- Zhou, X., Zhu, Y., Hou, D., Luo, Y., Peng, J. & Wobus, R. (2017) Performance of the NCEP global ensemble forecast system in a parallel experiment. *Weather Forecasting*, 32, 1989–2004.
- Zorita, E. & von Storch, H. (1999) The analog method as a simple statistical downscaling technique: comparison with more complicated methods. *Journal of Climate*, 12, 2474–2489.

How to cite this article: Feng, J., Toth, Z., Zhang, J. & Peña, M. (2024) Ensemble forecasting: A foray of dynamics into the realm of statistics. *Quarterly Journal of the Royal Meteorological Society*, 1–24. Available from: <https://doi.org/10.1002/qj.4745>

APPENDIX A. EXPERIMENTAL DATA

Experimental results in this study are based on operational analysis and forecast data from the NCEP Global Ensemble Forecast System (GEFS), initialized twice a day (0000UTC and 1200UTC) from the period December 1, 2017–February 28, 2018, for a total of 180 cases on a $1^\circ \times 1^\circ$ latitude–longitude grid, out to 16 days lead time at 12-hour output frequency (Zhou *et al.*, 2017). Note that the unperturbed control forecast is run at the same resolution as the perturbed forecasts. Most statistics are computed over the NH extratropics in the 30° – 65° latitude band. The perturbation methods and numerical model used to generate the NCEP ensemble are typical of those used at many other centers.

As reality (or truth) is unknown, true error cannot be measured in practice. In this study, we use NWP analysis fields as a proxy for truth. The difference between a forecast and this proxy can be called “perceived” error. With some assumptions, true error can be estimated based on perceived error measurements (Peña & Toth, 2014). Despite quantitative differences at short lead times, the qualitative behavior of true and perceived error is similar (Feng *et al.*, 2020). Beyond two days lead time, the bias in perceived forecast error induced by error in the verifying analysis field used as a proxy for truth is relatively small.

APPENDIX B. INFORMATION AND NOISE IN SIGNAL PROCESSING VERSUS WEATHER FORECASTING

Here we discuss what is common in and different between Information (I) as defined by Equation (7) and “information entropy” or Shannon entropy (SE, Shannon 1948) as used in information theory, and Noise as defined in Equation (8) compared with its use in signal processing. Both Information and SE provide a measure of uncertainty in our knowledge of a particular event out of all of its possible outcomes. While SE was introduced in the context of communication, Information is designed to quantify knowledge captured in analyzed or forecast states of a natural system like the atmosphere. Conveniently, I in its standardized form captures the fraction of forecast variance identical to the real state.

Noise, either defined by Equation (8) (N) or as used in signal processing, refers to impediments to accessing information. “In signal processing, noise is a general term for unwanted (and, in general, unknown) modifications that a signal may suffer during capture, storage, transmission, processing, or conversion” (Tuzlukov, 2010). Meanwhile, Noise in the context of forecast states of dynamical systems refers to dynamically constrained forecast variance that is unrelated to reality (see Section 4.3.4).

APPENDIX C. ERROR, NOISE, AND INFORMATION

Following Lorenz (1982), we assume that the divergence of initially nearby segments of a chaotic dynamical system’s trajectory, and in the absence of model error, true forecast error (i.e., the difference between a forecast and reality) follows a logistic curve:

$$d_i^2 = R \cdot c / (e^{-\alpha \cdot i \cdot \Delta t} + c), \quad (C1)$$

where $c = d_0^2 / (R - d_0^2)$, d_0^2 is the variance of initial error, R is the range between the lower and upper saturation values (that is double the climatic variance; Leith, 1974), α is the exponential growth rate, and t is the time increment.

Error variance (d_i^2) can also be expressed as a function of Information I_i , that is, the variance of truth missed by, and Noise variance (N_i) that is included in a forecast (see the top right-angled triangle in Figure 6):

$$d_i^2 = \frac{|\mathbf{F}_i - \mathbf{T}|^2}{|\mathbf{T} - \mathbf{C}|^2} = N_i + (1 - \sqrt{I_i})^2. \quad (C2)$$

For forecast systems with realistic variability, exploiting Equation (9), error variance can be written as a function of either Noise (not shown) or Information variance only:

$$d_i^2 = 2(1 - \sqrt{I_i}). \quad (C3)$$

Considering also Equations (C1) and (9), a rearrangement of Equation (C3) defines the time evolution of Noise (not shown) and Information (see blue line in Figure C1) as:

$$I_i = (2 - d_i^2)^2 / 4. \quad (C4)$$

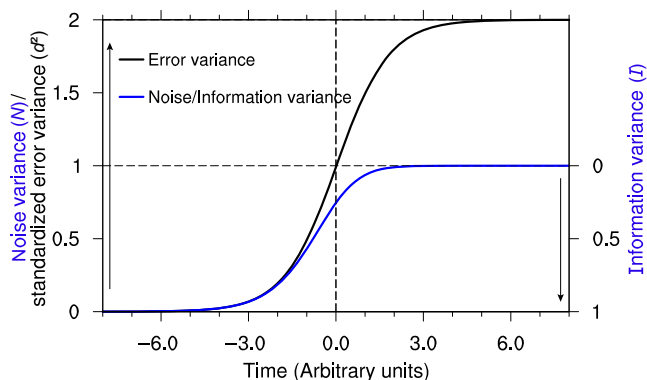


FIGURE C1 Schematic depicting the growth of Noise (blue line, left axis) and the decrease of Information variance (blue line, right axis) in a forecast characterized by logarithmically growing standardized error (black line). For further details, see text.

APPENDIX D. DEGREES OF FREEDOM

The experiment reported in Figure 9c is repeated with different values for dof and the frequency distribution of error amplitudes in perturbed states in the perfect (Figure 9b) and simulated ensembles are compared using the Kolmogorov–Smirnov two-sample test (Chakravarti *et al.*, 1967). Error amplitudes in both the perfect and simulated ensembles are standardized by the sample-mean rms error of the control analysis. The best fit is found at $M_d = 33$ (used in the construction of Figure 9c, see Figure D1a), with a range of values between 28 and 38 still acceptable at the 5% statistical significance level. To reduce noise, the test statistic is processed with a five-point triangular filter before it is plotted as a function of dof (Figure D1b). The results indicate that the experimental data in Figure 9b are consistent with the hypothesis that the global ensemble perturbations form a random sample in a high-dimensional phase space.

The $M_d = 33$ degrees of freedom (dof) estimated for the NH 500-hPa extratropical height, of course, assess only a small part of the full space of atmospheric dynamics at the resolution of today's models. Using the statistical evaluation described above, the best estimate for the dof of global 500-hPa height variability is found to be 50. Though global extratropical 500-hPa height covers a large subspace of atmospheric dynamics, it does not reflect independent variations across the entire planetary circulation (Palmer *et al.*, 2006, see their Appendix). Due to strong dynamical connections across variables, conservatively we expect an increase with a factor of less than 2 in dof if all independent model variables are considered. And due to the

low aspect ratio of the atmospheric fluid at today's resolution of global models (e.g., Held, 2015), we anticipate a similarly low (less than a factor of 2) increase in dof were all levels included. Such considerations suggest that the dof of the subspace of Information (i.e., initial error and short-range perturbation dynamics resolved by today's operational forecast systems) may be 3–4 times higher than that of the global 500-hPa height field, in the range of $M_d^{\text{overall}} = 150\text{--}200$.

APPENDIX E. BRACKETING RATIO IN MULTIPLE DIMENSIONS

Bracketing ratio F_{M_d, M_e} is a positively oriented metric, defined here for multidimensional applications as the relative frequency of reality (or its proxy) falling within (or bracketed by) the ensemble cloud in the direction congruent with the error in the control (see Appendix D). The bracketing ratio is a function of the degrees of freedom (M_d) and the number of ensemble members (M_e). Note that bracketing ratio F_{M_d, M_e} (Section 6.3) is an inverse measure of the ensemble outlier statistic (e.g., Buizza & Palmer, 1998), generalized for multidimensional applications, as well as a generalization of the probability of an ensemble member having an error lower than that in the control, shown in the table of the Appendix in Palmer *et al.* (2006).

For the illustration below, F_{M_d, M_e} is calculated as follows. Missed Information in the control and initial ensemble perturbation vectors d_0 and ε_0 is given by $(d_{0,1}, d_{0,2}, \dots, d_{0, M_d})$ and $(\varepsilon_{0,1}, \varepsilon_{0,2}, \dots, \varepsilon_{0, M_d})$, respectively. Since $\varepsilon_{0,i}$ is a random sample of $d_{0,i}$, and following the

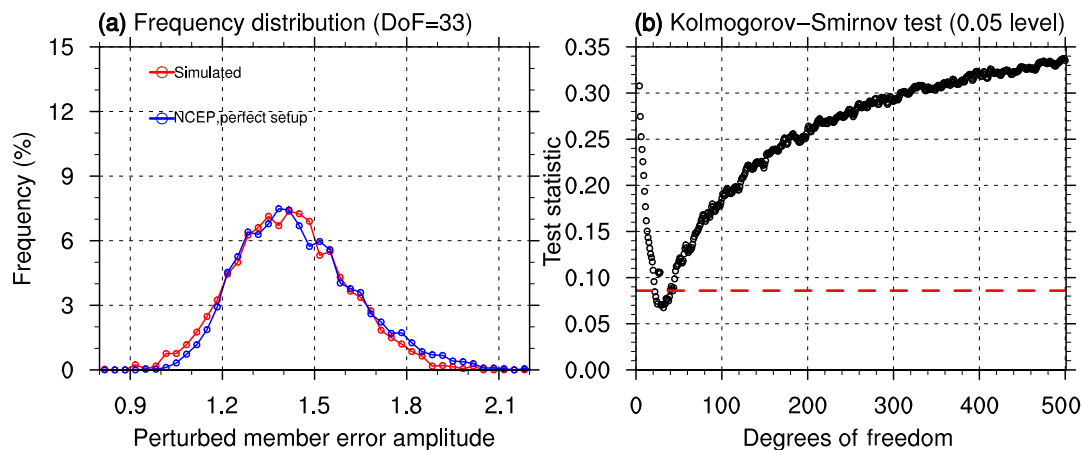


FIGURE D1 (a) The frequency of error in perturbed initial conditions from the NCEP (perfect setup, blue) and simulated (dof = 33, red) ensembles. (b) Test statistic for the two-sample Kolmogorov–Smirnov test showing the maximum absolute difference (black open circles) between the empirical perturbed-state error distribution functions from the perfect and simulated ensembles like those in panel (a). Values below the red dashed line indicate dof values where the actual and simulated distributions are statistically indistinguishable at the 0.05 significance level.

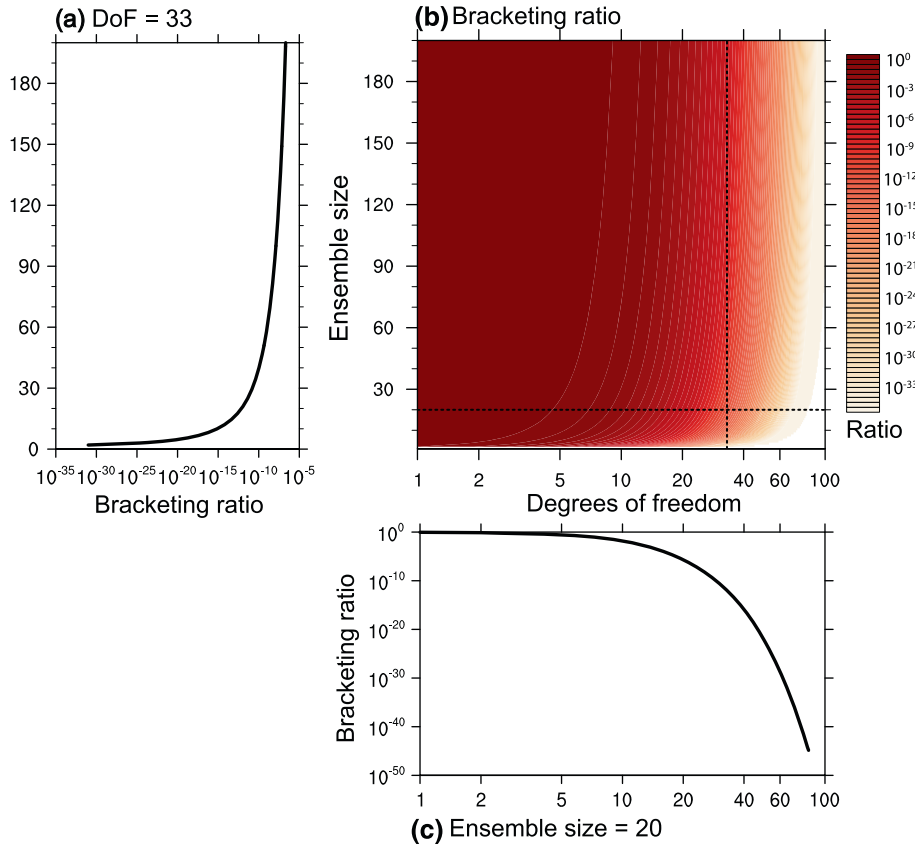


FIGURE E1 Ratio of cases a simulated ensemble of varying size brackets reality, as a function of the independent degrees of freedom (dof; panel b). Bracketing ratio for dof = 33 and a 20-member ensemble is highlighted in panels (a) and (c), respectively.

standardization introduced in Section 5.5, we assume the elements $\varepsilon_{0,i}$ and $d_{0,i}$ both follow independent and identical standard Gaussian distributions $N(0,1)$. Therefore, the distribution of projection of the ensemble perturbation ε_0 on the analysis error d_0 has an expected value of zero and a variance of 1, also conforming to a Gaussian distribution $N(0,1)$.

We consider an ensemble with M_e members. The projection of the members onto the direction congruent with the missed Information divide the probability space of $N(0,1)$ into $M_e + 1$ intervals. We mark the threshold designating the upper percentile of $1/(M_e + 1)$ as S . The truth is bracketed if the Euclidean norm of d_0 is smaller than S . The Euclidean norm of d_0 is calculated as $\sqrt{\sum_{i=1}^{M_d} d_{0,i}^2}$, where $\sum_{i=1}^{M_d} d_{0,i}^2$ follows the chi-squared distribution $\chi(M_d)$. Therefore, the general form of the formula for the bracketing ratio illustrated in Figure E1 is:

$$F_{M_d, M_e} = P(x < S^2), \quad (\text{E1})$$

where $x = \sum_{i=1}^{M_d} d_{0,i}^2 \sim \chi(M_d)$ and $P(\cdot)$ stands for the probability of x being smaller than S^2 . For a single variable

($M_d = 1$), Equation (E1) recovers the inverse of the formula for the often-used ensemble outlier statistic:

$$F_{1, M_e} = 1 - 2/(M_e + 1). \quad (\text{E2})$$

As an illustration, Figure E1b displays the expected value of F_{M_d, M_e} as a function of the degrees of freedom (M_d) and the number of ensemble members (M_e). Highlighted are marginal values for $M_d = 33$, the estimated dof of the NH extratropical 500-hPa height field (Figure E1a), and $M_e = 20$, the membership of the NCEP ensemble over the experimental period (Figure E1c). In sharp contrast with realistic-size ensembles in low dimensions ($F_{M_d, 20} \sim 1$, Figure E1c), even for large ensembles (e.g., $M_e = 200$) and for a limited domain like the NH extratropical 500-hPa height ($M_d = 33$), truth is bracketed only in one out of about 500 million cases (Figure E1a). This answers a question Gilmour and Smith (1997) posed in a broader context: ensembles can capture reality “only in” but not “even in” low-dimensional systems. In the full space of resolved atmospheric dynamics ($M_d^{\text{overall}} \sim 175$), truth would be encompassed at an astronomical rate so low that is computationally directly inaccessible.