



# Stochastic sampling provides a unifying account of visual working memory limits

Sebastian Schneegans<sup>a</sup>, Robert Taylor<sup>a</sup>, and Paul M. Bays<sup>a,1</sup>

<sup>a</sup>Department of Psychology, University of Cambridge, Cambridge CB2 3EB, United Kingdom

Edited by Wilson S. Geisler, The University of Texas at Austin, Austin, TX, and approved July 14, 2020 (received for review March 6, 2020)

Research into human working memory limits has been shaped by the competition between different formal models, with a central point of contention being whether internal representations are continuous or discrete. Here we describe a sampling approach derived from principles of neural coding as a framework to understand working memory limits. Reconceptualizing existing models in these terms reveals strong commonalities between seemingly opposing accounts, but also allows us to identify specific points of difference. We show that the discrete versus continuous nature of sampling is not critical to model fits, but that, instead, random variability in sample counts is the key to reproducing human performance in both single- and whole-report tasks. A probabilistic limit on the number of items successfully retrieved is an emergent property of stochastic sampling, requiring no explicit mechanism to enforce it. These findings resolve discrepancies between previous accounts and establish a unified computational framework for working memory that is compatible with neural principles.

visual working memory | population coding | resource model | capacity limits

Working memory refers to the nervous system's ability to form stable internal representations that can be actively manipulated in the pursuit of behavioral goals. A classical view of visual working memory (VWM) held that it was organized into a limited number of memory slots, each capable of holding a single object (1, 2). This model was subsequently modified to allow multiple slots to hold the same object and be combined on retrieval to achieve higher precision (3). This "slots+averaging" model incorporated aspects of an alternative view, which holds that VWM is better conceptualized as a continuous resource that can be flexibly distributed between different objects or visual elements (4, 5), accounting for set size effects in delayed reproduction tasks (6) (Fig. 1A) and flexibility in prioritizing representations (7). Variable precision models (8, 9) additionally proposed that the amount of memory resource is not fixed but varies randomly from item to item and trial to trial. An alternative approach (10) sought to explain VWM errors from neural principles as decoding variability in population representations (11), with the limited memory resource equated to the total neural activity dedicated to storage. Here we show that each of these influential accounts of VWM can be interpreted within a common framework based on the statistical principle of sampling (12–18).

## Sampling Interpretation of Population Coding

We first show how a population coding model (10) can, with some simplifying assumptions, be reinterpreted in terms of sampling (Fig. 1A–C). We consider a mathematically idealized population of independent neurons encoding a one-dimensional (1D) stimulus feature  $\theta$ , where the amplitude of each cell's activity is determined by its individual tuning function. Neurons are assumed to share the same tuning function, merely shifted so the peak lies at each neuron's preferred feature value  $\varphi_i$ ,

$$f_i(\theta) = f(\theta - \varphi_i). \quad [1]$$

Discrete spikes are generated from the cells' activity via independent Poisson processes. If we pick, at random, any spike generated by the neural population in response to a stimulus value  $\theta$ , we can determine the probability that it was produced by a neuron with preferred feature value  $\varphi$ . If we assume dense uniform coverage of the underlying feature space by neural tuning curves, this yields a continuous probability distribution  $p(\varphi)$  over the space of preferred feature values (Fig. 1C). This distribution has the same shape as the neural tuning curves and is centered on the true stimulus value,

$$p(\varphi) \propto f(\theta - \varphi). \quad [2]$$

Thus, if we associate each spike with the preferred feature value of the neuron that generated it (the principle of population vector decoding; ref. 19), we can interpret the spiking activity of the population as a set of noisy samples of the true stimulus value, drawn from the distribution  $p(\varphi)$ .

Retrieval of a feature value is modeled as decoding of the spikes generated within a fixed time window. In the idealized case with Gaussian tuning functions, the maximum likelihood (ML) decoder generates an estimate by simple averaging of the spike values,

$$\hat{\theta}_{ML} = \frac{1}{n} \sum_j^n \varphi_{(j)}, \quad [3]$$

where  $\varphi_{(j)}$  is the preferred feature value of the neuron that generated the  $j$ th spike.

Due to the superposition property of Poisson processes, the number of spikes—or samples—generated by the neural

### Significance

We demonstrate that three of the most prominent accounts of visual working memory in the psychology and neuroscience literature—the slots+averaging model, the variable precision model, and the population coding model—can all be expressed in the common mathematical framework of sampling. This reformulation allows us to pinpoint the key differences between these models, and to determine which factors are critical to account for the observed patterns of recall errors across different human psychophysical experiments. Moreover, the sampling framework provides a possible neural grounding for these models in the spiking activity of neuronal populations, as well as a link to existing theories of capacity limits in visual attention.

Author contributions: S.S., R.T., and P.M.B. performed research; S.S., R.T., and P.M.B. developed the theory; S.S. and P.M.B. carried out the analysis; S.S. and P.M.B. wrote the paper; and P.M.B. conceived and directed the project.

The authors declare no competing interest.

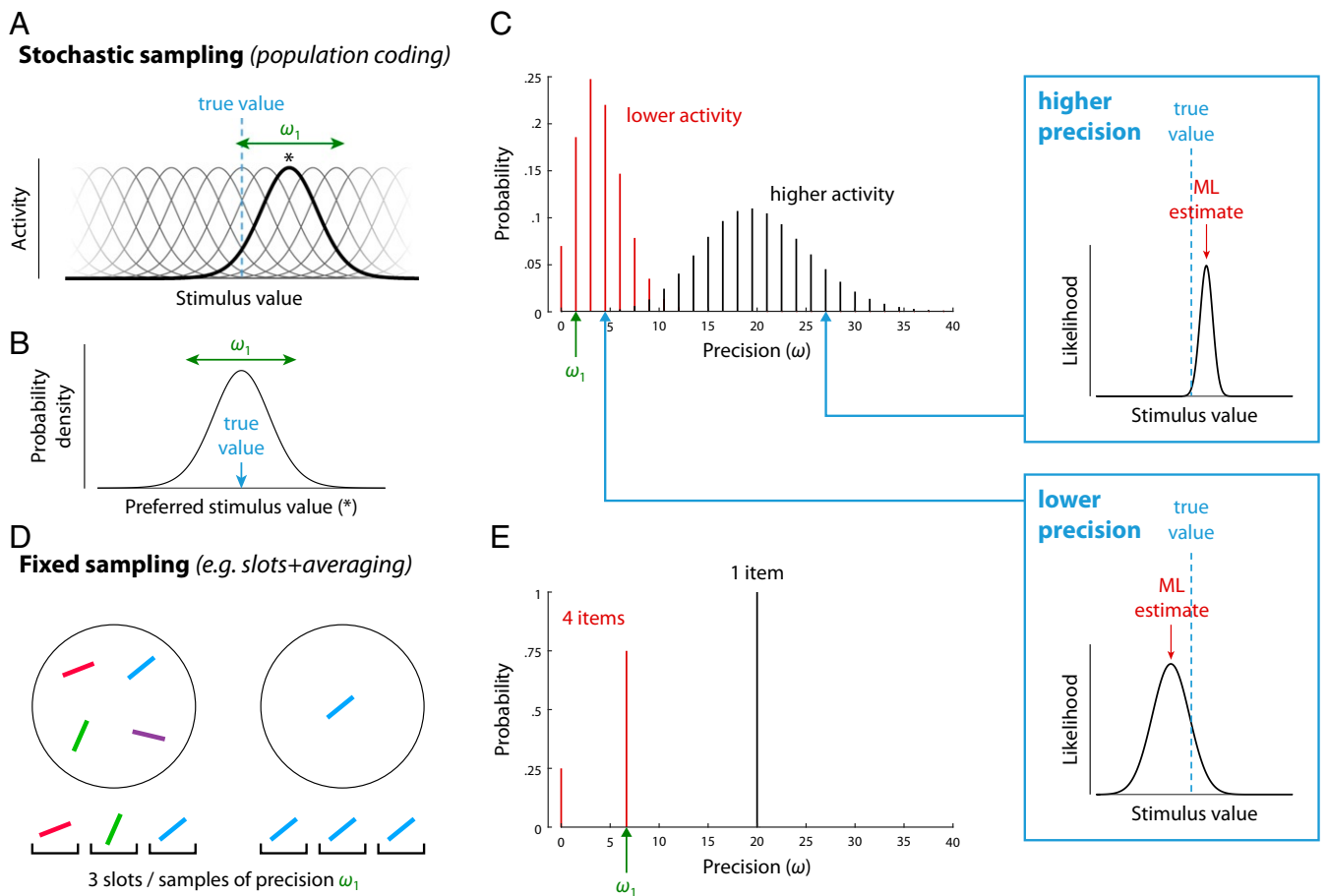
This article is a PNAS Direct Submission.

This open access article is distributed under Creative Commons Attribution License 4.0 (CC BY).

<sup>1</sup> To whom correspondence may be addressed. Email: pmb20@cam.ac.uk.

This article contains supporting information online at <https://www.pnas.org/lookup/suppl/doi:10.1073/pnas.2004306117/-DCSupplemental>.

First published August 11, 2020.



**Fig. 1.** Sampling interpretation of working memory models. (A–C) A theoretical account of neural population coding can be reinterpreted as sampling. (A) The stimulus-evoked response of spiking neurons in an idealized population depends on their individual tuning (one neuron’s tuning function and preferred value [\*] is highlighted). (B) Probability distribution over stimulus space obtained by associating a spike with the preferred stimulus of the neuron that generated it. (C) Precision of maximum likelihood estimates based on spikes emitted in a fixed decoding window. Precision, defined as the width of the likelihood function (*insets*), is discretely distributed as a product of the tuning precision ( $\omega_1$ ) and the number of spikes, which varies stochastically. Assuming normalization of total activity encoding multiple items, larger set sizes correspond to less mean activity per item. (D and E) An account based on averaging limited memory slots can also be described as sampling. (D) Allocation of a fixed number of samples or slots (here, three) to memory displays of different sizes. (E) Precision is discretely distributed as a product of the tuning width,  $\omega_1$ , and the number of samples allocated per item.

population within the decoding window is also a Poisson random variable. If the total spike rate in the neural population is normalized (20), or fixed at a population level  $\gamma$ , it implements a form of limited resource (10). This resource is continuous—unlike the discrete number of samples—and can be distributed between memory items, depending on task demands (e.g., prioritizing one item that is cued as a likely target). We will focus on the simplest case, in which the total spike rate is distributed evenly among all memory items, resulting in a mean number of samples available for decoding each stimulus that is inverse to the set size  $N$ . This has been shown to quantitatively capture the set size effect in single-report delayed reproduction tasks (Fig. 24).

The actual number of samples available in this model for decoding each item in a single trial,  $n_k$ , is a discrete random variable independently drawn from a Poisson distribution, with its mean determined by the spike rate for that item,

$$n_k \sim \text{Pois} \left( \frac{\gamma}{N} \right). \quad [4]$$

The neural population model can therefore be interpreted as a *stochastic* sampling model.

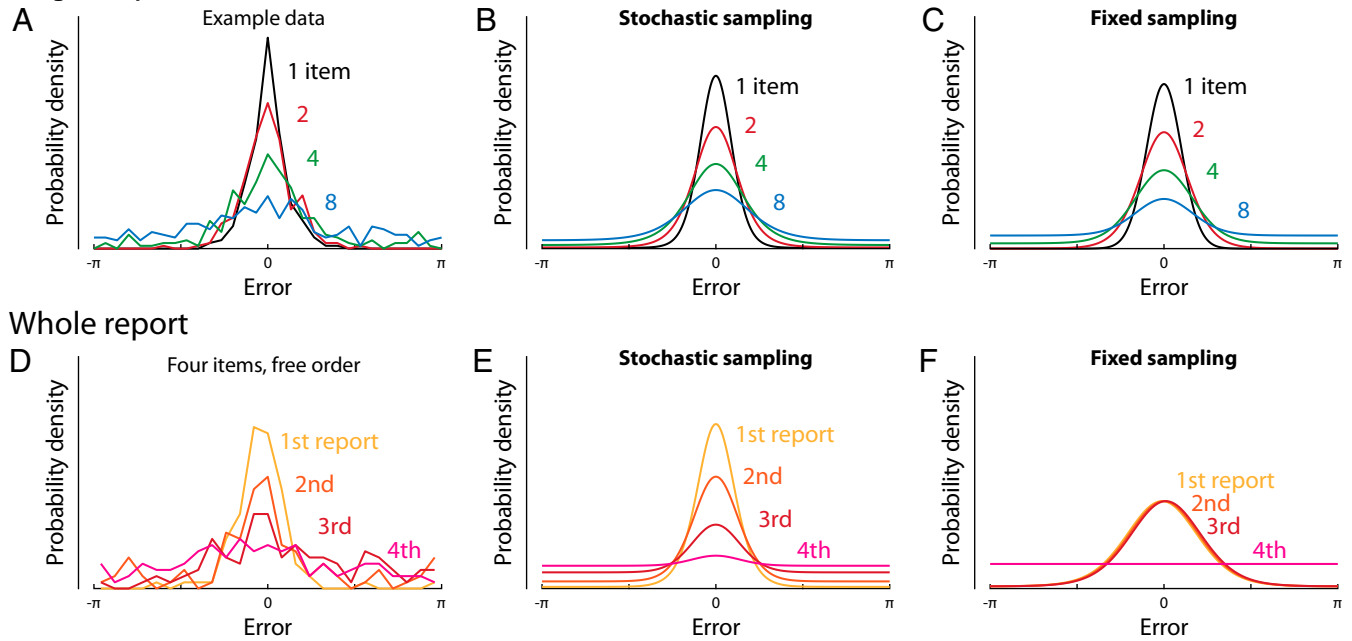
### Fixed Sampling Models

The most prominent discrete representation account of VWM, the slots+averaging model (3), can also readily be interpreted in terms of sampling (Fig. 1 D and E). Each slot is postulated to hold a representation of a single item with a fixed precision, and so provides a noisy sample of the item’s feature value (or values; the sampling interpretation is agnostic as to feature- vs. object-based views of VWM; refs. 21 and 22). Multiple slots, or samples, that correspond to the same object are averaged at retrieval to enhance the precision of the estimated stimulus feature. Thus, the format of representation and the decoding mechanism are identical to the stochastic sampling model. There is one critical difference, however: The slots+averaging model assumes that the total number of samples available for representing multiple items is fixed, that is,

$$\sum_k^N n_k = K. \quad [5]$$

This has also been the most common assumption in previous sampling-based models in the attentional and memory literature (refs. 12–14, but see ref. 23). We will refer to this as a *fixed* sampling model.

## Single Report



**Fig. 2.** Response distributions and model fits in delayed reproduction tasks. (A) Distributions of response errors in a single-report task for a representative participant at different set sizes (10). (B and C) ML fits of the data in A with the stochastic sampling model and fixed sampling model, respectively. (D) Distributions of response errors in a whole-report task for a representative participant at set size four, showing how errors increase with the (freely chosen) order of sequential report (24). (E and F) ML fits of the participant's data with the stochastic sampling model and fixed sampling model, respectively. Fits are based on results from all set sizes, not only the single set size shown in D.

### Predictions for Precision and Error Distributions

We now consider the distribution of representational precision in these models. For any particular set of samples,  $\phi$ , the information they provide about the stimulus is described by the likelihood function,  $\mathcal{L}(\theta; \phi) = p_{\theta}(\phi|\theta)$ , equivalent to the conditional probability of obtaining those samples given different values of the stimulus. The width of the likelihood function is a measure of uncertainty in the estimate: A set of samples with a broad likelihood function (Fig. 1 C, *Bottom Inset*) is compatible with many different feature values, whereas a narrow likelihood function (Fig. 1 C, *Top Inset*) identifies a value more precisely. While a pattern of samples may have a sharp likelihood function with a peak far from the true estimate (a kind of “false alarm”), statistically, this is unlikely.

If the sample values follow a normal distribution with variance  $\sigma^2$  centered on the true stimulus value, then the likelihood function is also normal, with a width that depends only on the number of samples available for decoding,

$$\mathcal{L}(\theta; \phi) \propto \phi \left( \theta; \hat{\theta}_{ML}, \frac{\sigma^2}{n_k} \right). \quad [6]$$

Furthermore, for a specified number of samples, the ML estimate is distributed around the true stimulus value as a normal with the same width as the likelihood,

$$\hat{\theta}_{ML}|n_k \sim \mathcal{N} \left( \theta, \frac{\sigma^2}{n_k} \right). \quad [7]$$

This correspondence between uncertainty, as expressed in the likelihood width, and trial-to-trial variability is not universal, but does apply to all of the models considered in this study, and justifies defining the precision of an individual estimate (which we will denote  $\omega$ ) as the precision of its corresponding likelihood function (see *SI Appendix, Fig. S1* for a detailed illustration).

Adopting this definition explicitly (see also refs. 25 and 26) allows us to treat precision as a random variable with a defined probability function, describing variation in the reliability of estimates while also predicting the distribution of errors across trials. This will prove critical in fitting data from whole-report tasks (Fig. 2D and below).

For the stochastic sampling model based on population coding, likelihood precision has a Poisson distribution (Fig. 1C), scaled by the precision of a single sample which is determined by the neural tuning function,  $\omega_1 = 1/\sigma^2$ ,

$$\frac{\omega}{\omega_1} \sim \text{Poisson} \left( \frac{\gamma}{N} \right). \quad [8]$$

Example distributions of decoding error are shown in Fig. 2B and E, where we have made a transition from 1D Euclidean to a circular stimulus space, corresponding more closely to the feature dimensions (e.g., orientation, hue) commonly used experimentally. The distribution of errors can be described as a scale mixture of normal distributions with precision proportional to the sample count (*SI Appendix, Fig. S1*; due to the circular stimulus space, this is a close approximation rather than exact; see *SI Appendix, Supporting Information Text*). The dispersion of errors increases with decreasing activity (e.g., as a result of increasing set size; Fig. 2B), and the distribution deviates from normality, with this effect being particularly evident at lower activity levels (blue curve) where long tails are observed.

For the fixed sampling model, making the common assumption that samples are distributed as evenly as possible among items (9, 27), we obtain a discrete distribution over, at most, two precision values (Fig. 1E), which are multiples of the precision of one sample,  $\omega_1$ . As in the stochastic model, mean precision is inversely proportional to set size, but, because the distributions over precision differ, the fixed and stochastic models make distinct, testable predictions for error distributions (Fig. 2).

## Response Errors Discriminate between Models

We tested the ability of stochastic and fixed sampling models to capture response errors in delayed reproduction tasks (*SI Appendix*, Fig. S2). We fit the models to a large dataset of single-report tasks originating from different laboratories (*SI Appendix*, Table S1) and also to a set of whole-report experiments (24), in which participants reported the feature values of all items in a sample array, either in a prescribed random order or in an order freely chosen by each participant on each trial. While only a single study, the whole-report results include information regarding correlations in errors between items represented simultaneously in VWM that could differentiate the models. On free choice trials, we assumed that participants gave their responses in order of decreasing precision (corresponding to decreasing number of samples and increasing likelihood width). This assumption is supported by previous findings that humans have knowledge about the uncertainty with which individual items are recalled (8, 25).

Overall, the stochastic model fit data substantially better than the fixed sampling model for both single-report (Fig. 3A; difference in log likelihood per participant,  $\Delta LL = 16.3 \pm 2.37$  [M  $\pm$  SE]) and whole-report tasks (Fig. 3B;  $\Delta LL = 162 \pm 13.6$ ), indicating that stochasticity is critical for capturing behavioral performance (see also *SI Appendix*, Fig. S3). The response error distributions in the whole-report task with freely chosen response order have previously been argued to provide evidence for a fixed item limit (24), since they approach uniform distributions for the later responses at high set sizes (Fig. 2D; see *SI Appendix*, Figs. S4 and S5 for full behavioral results and model fits). However, this qualitative observation is also predicted by the stochastic sampling model with responses ordered by precision, as the lowest precision retrievals will be based on few, or no, samples (Fig. 2E). Quality of fits could be further improved by taking into account deterioration of recall precision with increasing retention intervals (*SI Appendix*, Figs. S3J and S5), modeled as random drift of encoded feature values over time (28) (*SI Appendix*).

In contrast, the quantitative changes in error distribution with response order and set size were relatively poorly fit by the fixed sampling model (Fig. 2F). In particular, when the set size exceeds the fixed sample count, each item is represented by either one or zero samples, so this model cannot reproduce the graded decline in precision with response order that is also present in individual participants' data (and does not merely arise at the group level due to averaging across participants with different capacities).

We tested two intermediate model versions in order to further dissociate the specific aspects in which the fixed and stochastic sampling models differ, and determine the significance of each for capturing human performance. In the *random-fixed* model, the total number of samples was fixed but distributed randomly between items. This model provided an improved fit to data compared to the fixed model with even allocation (moderately for single-report,  $\Delta LL = 3.07 \pm 1.10$ ; strongly for whole-report,  $\Delta LL = 112 \pm 11.9$ ), but was still substantially worse than the stochastic model in both cases (single-report,  $\Delta LL = 13.2 \pm 2.24$ ; whole-report,  $\Delta LL = 50.4 \pm 7.03$ ). In the *even-stochastic* model, the total number of samples was a Poisson random variable, but the samples were distributed as evenly as possible between items. This model achieved a better fit to single-report data than the stochastic model with independent sample counts for each item ( $\Delta LL = 3.57 \pm 0.697$ ), but provided a much worse fit to whole-report data ( $\Delta LL = 21.4 \pm 4.12$ ). Combining evidence across all participants and tasks, the stochastic model with independent sample counts was strongly preferred over this and the other alternative models (total  $\Delta LL > 1,450$ ; Fig. 3C).

## Generalizing the Stochastic Model

For the models examined above, typical fitted parameters indicate that estimates are based on relatively small numbers of

samples (e.g., mean of  $\sim 13$  samples based on fits to single-report data). One result is that the precision of decoded estimates could take on only a limited set of possible values, and error distributions reflect a discrete mixture of distributions with different widths. From a neural perspective, while consistent with the remarkable fidelity with which single neurons' activity encodes visual stimuli (29, 30), such small sample counts nonetheless seem unlikely when interpreted as spike counts (see *Toward Biophysically Realistic Models*). To investigate whether discreteness and/or low numbers of samples are important for reproducing human performance, we therefore implemented a generalization of the stochastic model in which the number of samples was free to vary.

The distribution over precision values in the generalized stochastic model was obtained as a scaling of the negative binomial distribution,

$$\frac{\omega}{\omega_1 p} \sim \text{NegBin} \left( \frac{\gamma}{N} \frac{1}{1-p}, p \right). \quad [9]$$

This distribution has previously been proposed to model neural spiking activity (31), and it retains the characteristic relationship between mean and variability in the scaled Poisson distribution: The Fano factor (the ratio of variance to mean) is constant, equal to the value of a single sample,  $\text{Var}[\omega]/E[\omega] = \omega_1$ . This distinguishes the stochastic models from the fixed sampling model, where the Fano factor is at or close to zero (mean  $\sim 0.25$  of  $\omega_1$  based on ML parameters and typical set sizes) and varies in an idiosyncratic manner between set sizes, due to the varying combinatorial possibilities of allocating a fixed number of samples to a fixed number of items (Fig. 3D, purple).

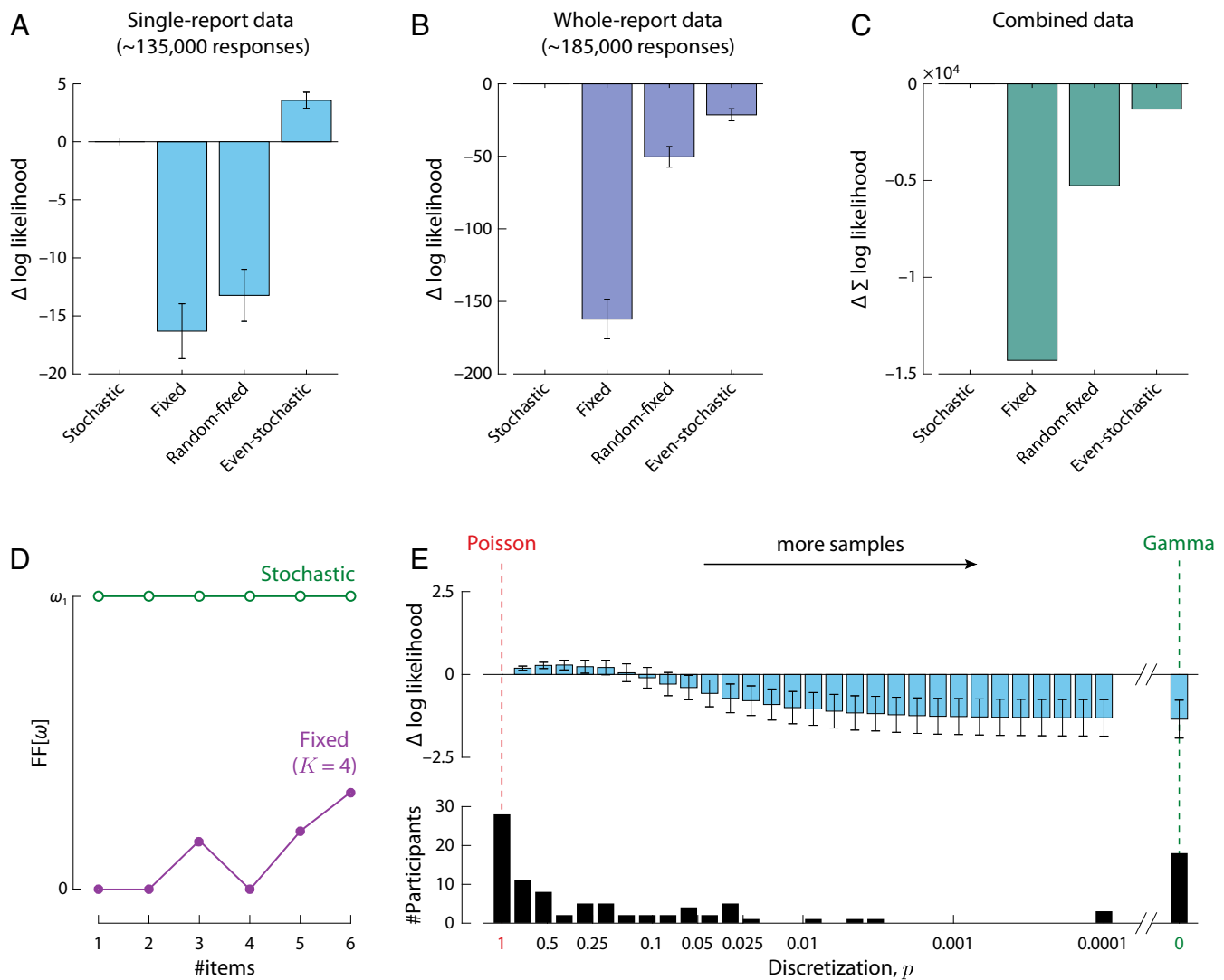
The parameter  $p$  in the generalized stochastic model controls the discretization of the precision distribution:  $p = 1$  corresponds to the Poisson model described above and illustrated in Fig. 4A (strictly, Eq. 8 is the limit of Eq. 9 as  $p \rightarrow 1$ ), while  $p < 1$  corresponds to a stochastic model with a greater mean number of samples,  $\bar{n} = \gamma/p$ , each with a lower individual precision,  $\omega_1 p$ . The mean and variance in precision ( $E[\omega] = \omega_1 \gamma/N$  and  $\text{Var}[\omega] = \omega_1^2 \gamma/N$ ), and thus also the Fano factor, are independent of the discretization  $p$ . Examples of precision distributions with different discretizations are shown in Fig. 4B and C.

As the discretization parameter becomes very small ( $p \rightarrow 0$ ), the number of samples becomes very large, and the distribution of precision described by Eq. 9 approaches a continuous function (Fig. 4D and *SI Appendix*), specifically the Gamma distribution,

$$\omega \sim \text{Gamma} \left( \frac{\gamma}{N}, \omega_1 \right). \quad [10]$$

Two previous studies (8, 9) independently proposed that a continuous scale mixture of normal distributions with Gamma-distributed precision provided a good account of VWM data, but did not provide a theoretical motivation for this choice of distribution. In particular, ref. 9 proposed distributing precision as  $\text{Gamma}(\bar{J}_1/N^\alpha, \tau)$ , with  $\bar{J}_1$ ,  $\tau$ , and  $\alpha$  as free parameters. With  $\alpha = 1$ , this is identical to Eq. 10 (see *SI Appendix* for results regarding this parameter). We can now explain Gamma-distributed precision as a limit case of the stochastic sampling model with large numbers of low-precision samples.

Fig. 3E, *Top* shows the results of fitting the generalized stochastic model with different levels of discretization,  $p$ , to the single-report dataset. The best fit was obtained with a discretization roughly one-third that of the Poisson model,  $p = 0.39$ . However, varying discretization produced differences in fit an order of magnitude smaller than those between fixed and stochastic sampling (varying by  $\sim 1.5$  versus  $\sim 15$  log likelihood points). Fitting the same model with  $p$  as a free parameter that could



**Fig. 3.** Model comparison based on single- and whole-report data. (A) Mean difference in log likelihood of each model from the stochastic sampling model (with independence between items), for a benchmark dataset of single-report experiments. More positive values indicate better fits to data. Error bars indicate  $\pm 1$  SE across participants. (B) The same comparison for a set of whole-report experiments. (C) Total difference in log likelihood between models across single- and whole report experiments. (D) Fano factor (ratio of variance to mean) of precision distribution. A constant Fano factor is characteristic of the stochastic model and contrasts with the varying Fano factor (dependent on set size and number of samples) in fixed sampling. (E) Mean difference in log likelihood for differing levels of discretization in the generalized stochastic model (Top), and number of participants best fit with each discretization level (Bottom). Differences in log likelihood are plotted relative to the maximum discretization ( $p = 1$ ; Left) corresponding to the standard stochastic model with Poisson-distributed precision. Lower discretization ( $p < 1$ ) corresponds to more samples each of lower precision, converging to a continuous Gamma distribution over precision as  $p$  approaches zero (Right). All models have the same number of free parameters and include a fixed per-item probability of swap errors (SI Appendix).

vary between participants, we found that ML estimates of discretization were very broadly distributed (Fig. 3 E, Bottom), with a majority of participants (72%) best described by a sampling model with less discreteness than the Poisson, and a minority (18%) better captured by the continuous limit ( $p \rightarrow 0$ ) than any discrete value of  $p$  we tested (as low as 0.0001, corresponding to  $\sim 100,000$  samples). Formal model comparison was equivocal with respect to an advantage of including the discretization parameter in comparison to either the Poisson model (i.e.,  $p = 1$ ; difference in Akaike Information Criterion,  $\Delta\text{AIC} = -0.61 \pm 0.49$ ; difference in Bayesian Information Criterion,  $\Delta\text{BIC} = +4.2 \pm 0.46$ ; negative values favor the added parameter) or the continuous Gamma model (i.e.,  $p \rightarrow 0$ ;  $\Delta\text{AIC} = -3.3 \pm 0.93$ ;  $\Delta\text{BIC} = +1.5 \pm 0.89$ ). Overall, these results do not allow strong conclusions to be drawn regarding

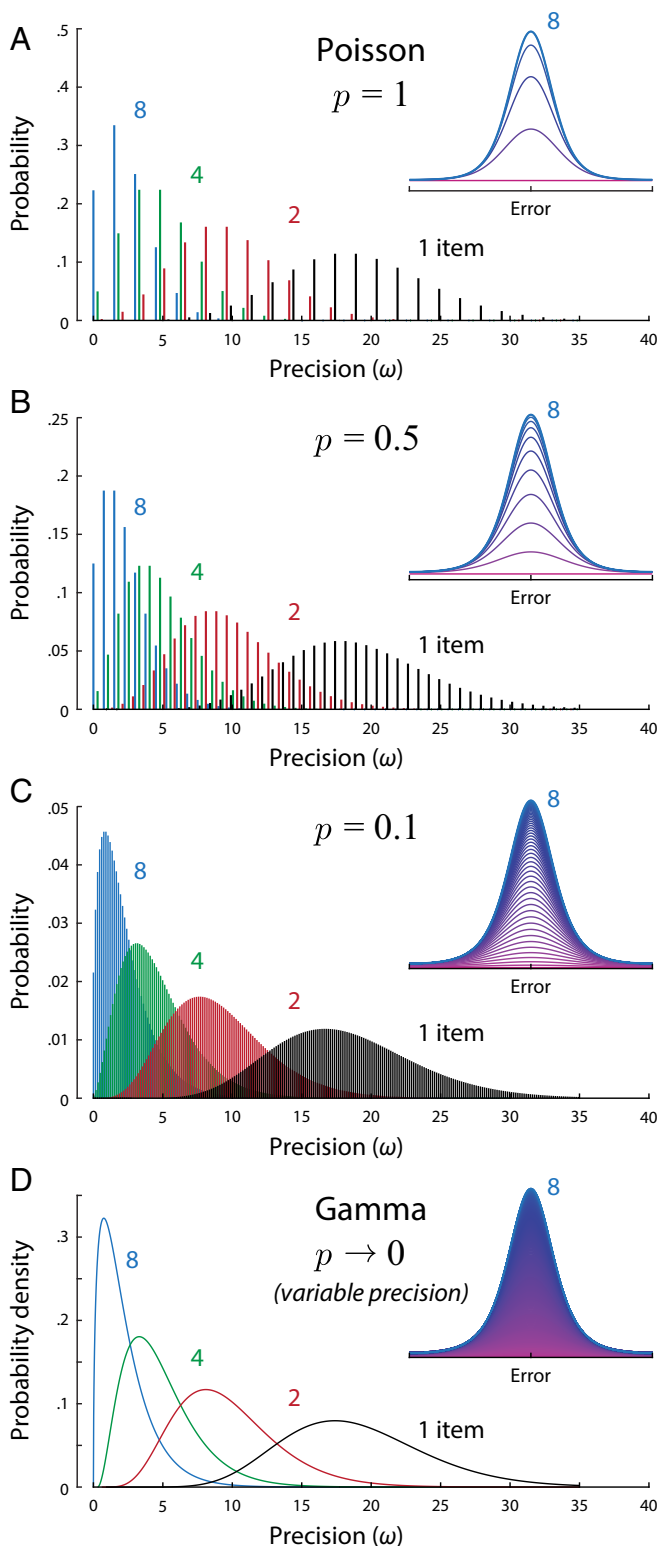
the discreteness of sampling, which has relatively little effect on error distributions (Fig 4, Insets) or the quality of fits.

### Probabilistic Item Limits

In the fixed sampling model, at higher set sizes, a meaningful proportion of estimates are random “guesses” based on no samples (Fig. 5 A and B). Specifically, if an estimate was generated for every item in the memory array, then, as set size  $N$  increased, the number of estimates based on at least one sample,  $S_{\omega>0}$ , would reach a maximum at the fixed total number of samples,

$$\lim_{N \rightarrow \infty} S_{\omega>0} = K. \quad [11]$$

This is a trivial consequence of sharing out a fixed number of samples evenly between items.



**Fig. 4.** Precision distributions in the generalized stochastic model, for different levels of discretization,  $p$ , and different set sizes. (Insets) Construction of the corresponding distributions of response error (for set size eight), with thin lines showing normal distributions with different precisions incrementally accumulated in ascending order (magenta to blue). (A) Example of discrete Poisson-distributed precision values ( $p = 1$ ). For typical ML parameters, estimates are based on a small mean number of samples (here,  $\gamma = 12$ ), each of moderate precision ( $\omega_1 = 1.5$ ). (B and C) With decreasing discretization ( $p < 1$ ), estimates are based on larger mean numbers of samples, and discrete precision values are more finely spaced. (D) In the limit as dis-

cretization falls to zero, the mean number of samples becomes infinite, and the distribution over precision approaches a continuous Gamma distribution. The ratio of variance to mean precision (Fano factor) is fixed (at  $\omega_1 = 1.5$ ) across all set sizes and levels of discretization.

$$\Pr(\omega = 0) = \exp\left(-\frac{\gamma}{N}\right), \quad [12]$$

and the number of nonrandom retrievals in a set of  $N$  items has a binomial distribution,

$$S_{\omega > 0} \sim \text{Bin}\left(N, 1 - \exp\left(-\frac{\gamma}{N}\right)\right). \quad [13]$$

As set size increases, the mean number of estimates based on at least one sample reaches a maximum at the expected total number of samples (Fig. 5 C and D). However, unlike the fixed sampling model, this limit is probabilistic, and (as illustrated in Fig. 5C) the actual number will vary from one set of memory items to the next, converging to a Poisson distribution for large  $N$ ,

$$\lim_{N \rightarrow \infty} S_{\omega > 0} \sim \text{Poisson}(\gamma). \quad [14]$$

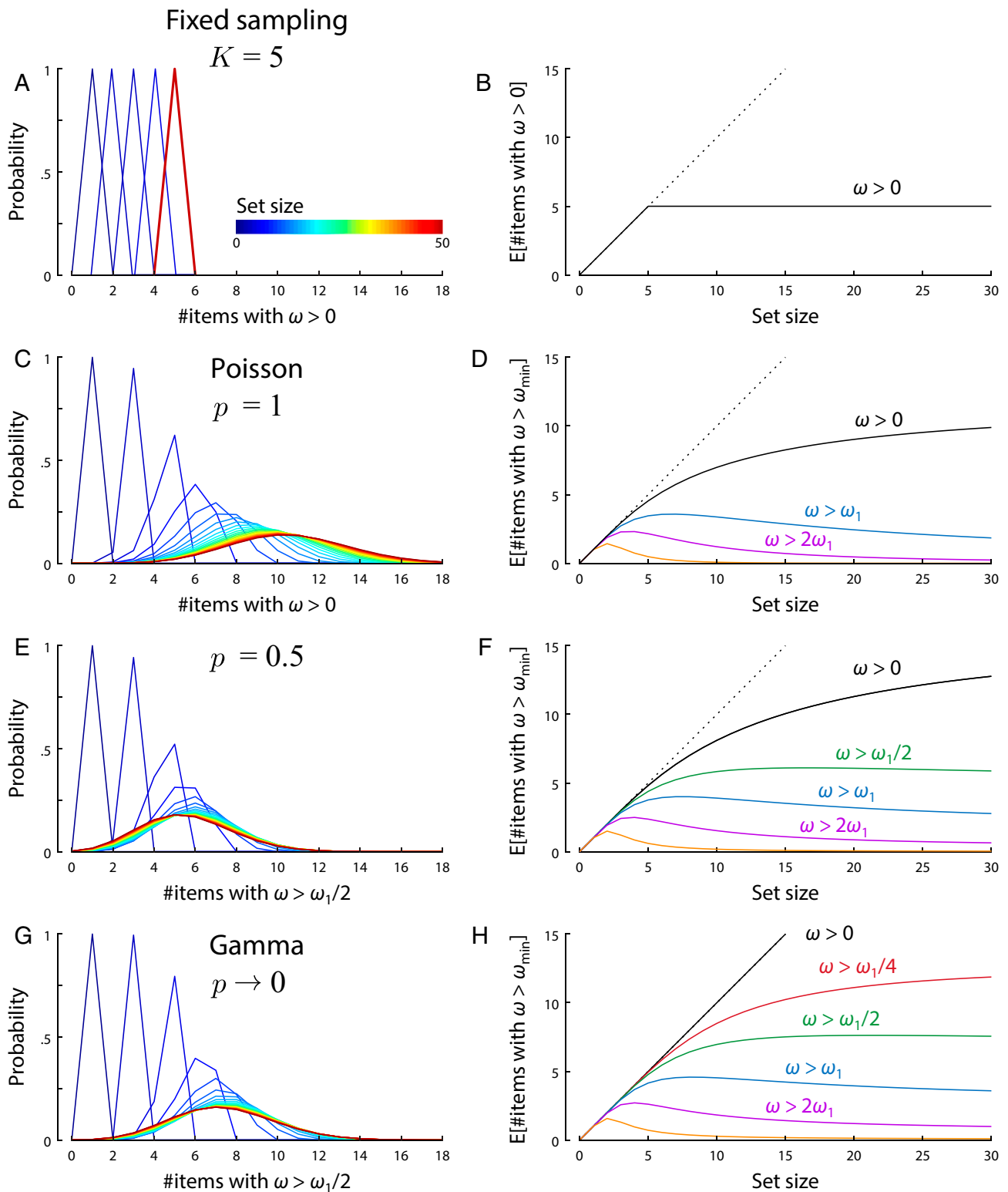
As we increase the expected number of samples by reducing the discretization,  $p$ , the probability of zero samples falls to zero:  $\Pr(\omega = 0) = p^{\gamma/(N(1-p))}$ . However, if we choose a precision threshold that is less than or equal to the base precision  $\omega_1$ , it can be shown that the mean number of items with above-threshold precision converges to a finite positive number at large set sizes (SI Appendix). This saturation is illustrated for different levels of discretization and various precision thresholds in Fig. 5 E–H.

Item limits or “magic numbers” (2, 24, 32) are usually considered synonymous with slot-based accounts, occurring when some items must go unrepresented because other items have filled the available capacity. The present results show that a probabilistic item limit, that is, an upper limit on the average number of items successfully retrieved that is not exceeded at any set size, can arise even when the probability of success for one item is independent of each other item. This holds true in the Poisson sampling model if we define success as obtaining one or more samples, but also more generally, even in a continuous model, if we define success as exceeding a threshold level of precision in estimation. Note, however, that the item limit does not, in general, have a simple relationship to the underlying number of samples. For example, the probabilistic limit at approximately five items in Fig. 5E is obtained from a model with a mean of 24 samples.

### Toward Biophysically Realistic Models

The idealized description of population coding on which we based the stochastic sampling model overlooks a number of important considerations in order to reveal relationships between cognitive and neural-level accounts of VWM. For instance, the statistics of spike counts in the neural system often deviate from the Poisson distribution assumed in the original population coding model in that they are “overdispersed” (i.e., Fano factor of  $> 1$ ). Such an overdispersion in the sample count also occurs in the generalized stochastic model as the discretization  $p$  decreases (in order for the Fano factor of the precision distribution to remain constant, the Fano factor of the sample count has to increase). Spike counts in visual cortical neurons typically show a Fano factor in the range 1.5 to 3 (e.g., ref. 33),

cretization falls to zero, the mean number of samples becomes infinite, and the distribution over precision approaches a continuous Gamma distribution. The ratio of variance to mean precision (Fano factor) is fixed (at  $\omega_1 = 1.5$ ) across all set sizes and levels of discretization.



**Fig. 5.** Item limits in sampling models. For each model, *A*, *C*, *E*, and *G* show how the probability distribution of the number of items recovered with greater than zero precision (*A* and *C*; greater than a fixed threshold for *E* and *G*) changes with set size (color coded, increasing blue to red; discrete probability distributions are depicted as line plots for better visualization). *B*, *D*, *F*, and *H* plot the mean number of items with above-threshold precision as a function of set size for different threshold values. Thresholds are defined as a proportion of the base precision  $\omega_1$ . (*A* and *B*) In the fixed sampling model, the number of items with nonzero precision increases with set size, then plateaus when the number of items equals the number of samples. (*C* and *D*) The stochastic sampling model with Poisson variability also has a limit on the number of items with nonzero precision, although this limit is probabilistic and emerges asymptotically (converging to the distribution shown by the red curve in *C* for large set sizes, corresponding to the mean number of items plotted as black curve in *D*). (*E* and *F*) Stochastic models with lower discretization display similar probabilistic item limits for precision exceeding a fixed threshold, but with the expected number of items saturating at different values depending on threshold (different colors in *F*). (*G* and *H*) This property also extends to models with continuous precision distributions.

corresponding, in our model, to discretization  $p$  in the range 0.33 to 0.75.

In real neural populations, there is also considerable variability between individual cells' tuning curves (34). Due to this heterogeneity, neurons differ in the amount of information each spike provides about a stimulus. From a sampling perspective, this means that estimates are based on samples that vary in precision, and this has the effect of "smoothing out" the discrete distribution of precision values predicted by the stochastic model (*SI Appendix*, Fig. S6). This has similar consequences for estimation error as decreasing  $p$  in the generalized model. We fit the single report data with a variant of the population model with random variability in the neurons' tuning curves (affecting baseline activity, gain, and tuning curve width, as well as adding heterogeneity in the coverage of the feature space by neural tuning curves; *SI Appendix*), scaled by a global heterogeneity parameter  $\nu$ . Incorporating biologically realistic heterogeneity into the population model improved fits to data ( $\Delta\text{AIC} = 8.3 \pm 1.8$ ,  $\Delta\text{BIC} = 3.4 \pm 1.7$  compared to the stochastic sampling model). The mean heterogeneity parameter in the ML fits was  $\nu = 0.66 \pm 0.08$ , where  $\nu = 0$  means no heterogeneity, and  $\nu = 1$  was approximately matched to heterogeneity of orientation-selective neurons in recordings from primary visual cortex (34).

Finally, spikes in real neural populations are not independent events as assumed by the sampling interpretation, but rather are correlated within and between neurons. This will tend to result in deviations from the simple additivity assumed by sampling. An implementation of short-range pair-wise correlations in the heterogeneous population model (see *SI Appendix* for details) greatly increased the numbers of decoded spikes required to reproduce behavioral data (on average, 168 times higher), without changing quality of fit ( $\Delta\text{AIC/BIC} = 0.045 \pm 0.28$ ). We note, however, that the exact consequences of spike correlations for decoding depend on details of correlation structure that are difficult to measure experimentally (35–37), and suboptimal inference (in the form of a mismatched decoder) may play an important part (38).

## Discussion

Taking, as a starting point, a mathematical idealization of the way neural populations encode information, we have shown that retrieval of a visual feature from working memory can be described as estimation based on a stochastically varying number of noisy samples. Two other influential models of VWM can be reconceptualized in the same framework: The slots+averaging model, because it modified the original slot model to allow multiple representations with independent noise, is directly equivalent to a sampling model with a fixed number of samples (13). And the variable precision model (8, 9) constitutes the continuous limit of a sampling model as samples are made less precise and more numerous, while maintaining the fixed proportionality between the variance and mean of precision in the decoded estimate.

Formulating all three models in the same mathematical framework (a formal "unification") allowed us to pinpoint specific differences between them. We determined the effect of these differences on the models' ability to account for human behavior by fitting multiple variants of the sampling model to a large database of delayed reproduction tasks. We found that stochasticity both in the total number of samples and in their distribution among items has a major impact on the quality of fit, with the best fits obtained if the number of samples is drawn randomly and independently for each item in each trial. Note that this form of stochasticity is poorly captured by the concept of memory "slots," because of the implication that a slot occupied by one item leaves fewer slots available for other items—this would predict dependencies between items in whole report that were not supported experimentally.

On the other hand, contrary to the assumptions of continuous resource models, we did find limited support for discreteness of memory representations (3). The fully continuous model with Gamma-distributed precision proposed in previous studies provided fits to data that were, overall, a little worse than the discrete Poisson model, in both single- and whole-report tasks (*SI Appendix*). When we attempted to fit discretization as a free parameter, however, we found that ML estimates varied widely between participants, and many were best fit by continuous or near-continuous versions of the generalized stochastic sampling model. So, while discreteness in memory representations is plausible—even inevitable if based on discrete spiking activity—recall errors do not provide strong evidence for any one particular level of discreteness, or, as a corollary, any particular mean number of samples.

Our findings further highlight the need to distinguish between two concepts that have previously been elided: discreteness in representation and discreteness in allocation. In the stochastic sampling model, the resource underlying capacity limits in VWM is equated with the mean number of samples (or the mean spike rate in the population coding interpretation), which can be distributed among items in a continuous fashion, even though the consequent number of samples obtained by each item is a discrete integer. This view on memory resources was strongly motivated by studies showing that prioritized items can be represented more precisely in VWM, at the cost of decreased precision for other items (7, 39–42). The stochastic sampling model can account for such findings through an uneven distribution of resources among memory items, corresponding to a higher mean number of samples for some items at the cost of a lower mean for other items. The actual number of samples available on an individual memory retrieval varies randomly about the item's mean. In the neural population model, this mechanism has previously been shown to successfully reproduce data from tasks in which one item is cued as the likely target (10).

While the stochastic sampling model is based on a highly idealized implementation of population coding, it nevertheless provides a link to a concrete neural mechanism that could form the basis of VWM performance. We have shown that adapting this model to achieve a higher degree of biophysical realism—by introducing heterogeneity in neural tuning curves and correlated spiking activity—improved the quality of fit to behavioral data. It has recently been shown that more neurally realistic population coding models preserve the key characteristics of the idealized model, and that signatures of neural tuning may even be visible in behavioral data (43).

Our results also provide a link between models of working memory used in the psychological literature and more biophysically detailed neural models such as continuous attractor networks (44–46), whose greater complexity typically precludes quantitative fits to behavioral data. These models are, likewise, based on principles of population coding and emphasize the role of neural noise in explaining variability in working memory performance. They are capable of producing probabilistic item limits similar to those described here, but it remains unresolved how these models could account for the graded variations in recall fidelity that we have found to be essential for capturing human behavioral performance.

In keeping with most previous work on VWM limits, we have not here attempted to reproduce the variations in bias and precision that are observed for different feature values, exemplified by the finding that cardinal orientations can be reproduced with greater precision than obliques. However, previous work has shown that these effects can be simply and elegantly captured within the population coding framework via the principle of efficient coding (47–49). The idea is that neural tuning functions are adapted to the stimulus statistics of the natural environment in such a way as to maximally convey information in that



environment (effectively by distributing neural resources preferentially to the most frequently occurring stimulus features). Although it should be possible to formulate this model as a modification of stochastic sampling, without reference to neural populations, it seems that the modifications required would not have a natural explanation within the sampling framework. These observations, and the results of incorporating heterogeneity described above, illustrate the value of connecting abstract cognitive models to neural theory.

We also did not address here the question of how individual features of a visual stimulus are bound together, which forms another point of contention in the debate on the format of VWM representations. In the model fits, we allowed failures of binding memory in the form of swap errors to occur with a fixed rate, although taking into account similarity of items with respect to the cue feature is likely to improve model fits (50, 51). While the discrete memory representations in slot models have traditionally been associated with a strongly object-based view (1), the sampling framework is agnostic as to whether objects or features are the units of VWM storage. Both views are compatible with the population coding interpretation, depending on whether the neurons in question are sensitive to a single feature (52) or a conjunction of features (51, 53).

A recent proposal that VWM errors can be explained in terms of a perceptual rescaling of stimulus space can also be expressed in terms of population coding, with some minor differences from the version presented here (see ref. 54 for details and discussion). In particular, the idea of retrieval based on normally distributed “memory-match” signals maps exactly onto an idealized population code with continuous-valued activity and constant Gaussian noise (51, 55). This predicts a continuous distribution over precision, not dissimilar to the Gamma distribution. Continuousness in representation does not appear a necessary component of this account, however, and it should be possible to reformulate it with arbitrary levels of discreteness, as in our generalized stochastic model.

There are other models of working memory that address capacity limits without explicitly postulating a limited memory resource (56). Some accounts stress the importance of memory decay over time, and active rehearsal to counteract this decay (57, 58). These theories do not have a clear analogue in the sampling framework, although effects of retention time have been incorporated into the neural population model (28). Other accounts have sought to explain capacity limits by interference between different memorized items (59). While the sampling framework does not explicitly address interference, the effect of normalization could be described as a form of nonspecific interference between items. A model of feature binding based on the neural population model shows some notable congruencies with

an interference account of VWM, and both models make similar predictions regarding swap errors (50, 51). Further research will be needed to determine the exact relationship between these models.

Taken together, our results reveal a surprising convergence between prominent models of VWM. Despite the fact that these competing models were independently motivated by different behavioral and neural findings, they can be expressed within the shared formal framework of sampling, which reveals specific distinguishing factors as well as shared general principles. This convergence gives cause for confidence that the stochastic sampling model captures key characteristics of VWM and will provide a solid foundation for future research.

## Materials and Methods

We fit computational models of VWM to behavioral data from a large dataset of delayed estimation experiments. The dataset included 15 individual single-report experiments (*SI Appendix, Table S1*; see *SI Appendix, Supporting Information Text* for inclusion criteria), as well as four whole-report experiments (*SI Appendix, Table S2*). Each model defines a parameterized distribution of response probabilities given the true feature values of the target and nontarget items in each trial (*SI Appendix*). For fits to whole-report data, we determined the probabilities of obtaining the given combination of responses within a single trial, taking into account the correlations of recall precision between different items within a trial predicted by each model.

We obtained an ML fit of each subject’s data for each model. The stochastic sampling, fixed sampling, and continuous sampling (Gamma) model, as well as the fixed-random and stochastic-even variants, each have three free parameters (including one parameter for the probability of swap errors). We fit these to both the single-report and whole-report data using the Nelder-Mead simplex algorithm (see *SI Appendix* for details). The generalized stochastic model and the neural population model with heterogeneous tuning curves have four free parameters each. For these models, we employed a grid search to obtain fits only of the single-report data (fitting them to whole-report data was not computationally feasible). We further evaluated model variants employing a more accurate method for ML decoding for circular feature spaces (rather than the Gaussian approximation used for fits reported in *Response Errors Discriminate between Models*), models without swap errors, models with an additional free parameter for the power law in set size effects, and models with a temporal decay of memory precision over varying response delays in the whole-report experiments (*SI Appendix*).

**Data Availability.** Data and code associated with this study are publicly available in Open Science Framework (60).

**ACKNOWLEDGMENTS.** We thank Ronald van den Berg, Wei Ji Ma, Máté Lengyel, and Masud Husain for helpful conversations, Zakhar Kabluchko and Martin Bays for statistical advice, and all of the researchers who publicly shared data that facilitated this study. We used resources provided by the Cambridge Service for Data Driven Discovery operated by the University of Cambridge Research Computing Service. This research was supported by the Wellcome Trust (Grant 106926).

- S. J. Luck, E. K. Vogel, The capacity of visual working memory for features and conjunctions. *Nature* **390**, 279–281 (1997).
- N. Cowan, The magical number 4 in short-term memory: A reconsideration of mental storage capacity. *Behav. Brain Sci.* **24**, 87–114 (2001).
- W. Zhang, S. J. Luck, Discrete fixed-resolution representations in visual working memory. *Nature* **453**, 233–235 (2008).
- P. M. Bays, R. F. G. Catalao, M. Husain, The precision of visual working memory is set by allocation of a shared resource. *J. Vis.* **9**, 7–7 (2009).
- W. J. Ma, M. Husain, P. M. Bays, Changing concepts of working memory. *Nat. Neurosci.* **17**, 347–356 (2014).
- P. Wilken, W. J. Ma, A detection theory account of change detection. *J. Vis.* **4**, 11 (2004).
- P. M. Bays, M. Husain, Dynamic shifts of limited working memory resources in human vision. *Science* **321**, 851–854 (2008).
- D. Fougny, J. W. Suchow, G. A. Alvarez, Variability in the quality of visual working memory. *Nat. Commun.* **3**, 1229 (2012).
- R. van den Berg, H. Shin, W.-C. Chou, R. George, W. J. Ma, Variability in encoding precision accounts for visual short-term memory limitations. *Proc. Natl. Acad. Sci. U.S.A.* **109**, 8780–8785 (2012).
- P. M. Bays, Noise in neural populations accounts for errors in working memory. *J. Neurosci.* **34**, 3632–3645 (2014).
- A. Pouget, P. Dayan, R. Zemel, Information processing with population codes. *Nat. Rev. Neurosci.* **1**, 125–132 (2000).
- M. L. Shaw, Identifying attentional and decision-making components in information processing. *Atten. Perform.* **8**, 277–295 (1980).
- J. Palmer, Attentional limits on the perception and memory of visual information. *J. Exp. Psychol. Hum. Percept. Perform.* **16**, 332–350 (1990).
- D. K. Sewell, S. D. Lilburn, P. L. Smith, An information capacity limitation of visual short-term memory. *J. Exp. Psychol. Hum. Percept. Perform.* **40**, 2214–2242 (2014).
- A.-M. Bonnel, J. Miller, Attentional effects on concurrent psychophysical discriminations: Investigations of a sample-size model. *Percept. Psychophys.* **55**, 162–179 (1994).
- E. Vul, N. Goodman, T. L. Griffiths, J. B. Tenenbaum, One and done? Optimal decisions from very few samples. *Cognit. Sci.* **38**, 599–637 (2014).
- G. Orbán, P. Berkes, J. Fiser, M. Lengyel, Neural variability and sampling-based probabilistic representations in the visual cortex. *Neuron* **92**, 530–543 (2016).
- M. N. Shadlen, D. Shohamy, Decision making and sequential sampling from memory. *Neuron* **90**, 927–939 (2016).
- A. P. Georgopoulos, A. B. Schwartz, R. E. Kettner, Neuronal population coding of movement direction. *Science* **233**, 1416–1419 (1986).
- M. Carandini, D. J. Heeger, Normalization as a canonical neural computation. *Nat. Rev. Neurosci.* **13**, 51–62 (2012).

21. G. A. Alvarez, P. Cavanagh, The capacity of visual short-term memory is set both by visual information load and by number of objects. *Psychol. Sci.* **15**, 106–111 (2004).
22. K. Oberauer, S. Eichenberger, Visual working memory declines when more features must be remembered for each object. *Mem. Cognit.* **41**, 1212–1227 (2013).
23. P. L. Smith, The Poisson shot noise model of visual short-term memory and choice response time: Normalized coding by neural population size. *J. Math. Psychol.* **66**, 41–52 (2015).
24. K. C. S. Adam, E. K. Vogel, E. Awh, Clear evidence for item limits in visual working memory. *Cognit. Psychol.* **97**, 79–97 (2017).
25. R. Van den Berg, A. H. Yoo, W. J. Ma, Fechner's law in metacognition: A quantitative model of visual working memory confidence. *Psychol. Rev.* **124**, 197–214 (2017).
26. P. M. Bays, A signature of neural coding at human perceptual limits. *J. Vis.* **16**, 4 (2016).
27. R. van den Berg, E. Awh, W. J. Ma, Factorial comparison of working memory models. *Psychol. Rev.* **121**, 124–149 (2014).
28. S. Schneegans, P. M. Bays, Drift in neural population activity causes working memory to deteriorate over time. *J. Neurosci.* **38**, 4859–4869 (2018).
29. E. Zohary, M. N. Shadlen, W. T. Newsome, Correlated neuronal discharge rate and its implications for psychophysical performance. *Nature* **370**, 140–143 (1994).
30. K. H. Britten, M. N. Shadlen, W. T. Newsome, J. A. Movshon, The analysis of visual motion: A comparison of neuronal and psychophysical performance. *J. Neurosci.* **12**, 4745–4765 (1992).
31. R. L. T. Goris, J. A. Movshon, E. P. Simoncelli, Partitioning neuronal variability. *Nat. Neurosci.* **17**, 858–865 (2014).
32. G. A. Miller, The magical number seven, plus or minus two: Some limits on our capacity for processing information. *Psychol. Rev.* **101**, 343–352 (1994).
33. R. Vogels, W. Spileers, G. A. Orban, The response variability of striate cortical neurons in the behaving monkey. *Exp. Brain Res.* **77**, 432–436 (1989).
34. A. S. Ecker *et al.*, Decorrelated neuronal firing in cortical microcircuits. *Science* **327**, 584–587 (2010).
35. B. B. Averbeck, P. E. Latham, A. Pouget, Neural correlations, population coding and computation. *Nat. Rev. Neurosci.* **7**, 358–366 (2006).
36. A. S. Ecker, P. Berens, A. S. Tolias, M. Bethge, The effect of noise correlations in populations of diversely tuned neurons. *J. Neurosci.* **31**, 14272–14283 (2011).
37. R. Moreno-Bote *et al.*, Information-limiting correlations. *Nat. Neurosci.* **17**, 1410–1417 (2014).
38. J. M. Beck, W. J. Ma, X. Pitkow, P. E. Latham, A. Pouget, Not noisy, just wrong: The role of suboptimal inference in behavioral variability. *Neuron* **74**, 30–39 (2012).
39. P. M. Bays, N. Gorgoraptis, N. Wee, L. Marshall, M. Husain, Temporal dynamics of encoding, storage, and reallocation of visual working memory. *J. Vis.* **11**, 6–6 (2011).
40. A. H. Yoo, Z. Klyszajko, C. E. Curtis, W. J. Ma, Strategic allocation of working memory resource. *Sci. Rep.* **8**, 16162 (2018).
41. S. M. Emrich, H. A. Lockhart, N. Al-Aidroos, Attention mediates the flexible allocation of visual working memory resources. *J. Exp. Psychol. Hum. Percept. Perform.* **43**, 1454–1465 (2017).
42. Z. Klyszajko, M. Rahmati, C. E. Curtis, Attentional priority determines working memory precision. *Vis. Res.* **105**, 70–76 (2014).
43. R. Taylor, P. M. Bays, Theory of neural coding predicts an upper bound on estimates of memory variability. *Psychol. Rev.*, 10.1037/rev0000189 (2020).
44. J. S. Johnson, J. P. Spencer, S. J. Luck, G. Schöner, A dynamic neural field model of visual working memory and change detection. *Psychol. Sci.* **20**, 568–577 (2009).
45. Z. Wei, X.-J. Wang, D.-H. Wang, From distributed resources to limited slots in multiple-item working memory: A spiking network model with normalization. *J. Neurosci.* **32**, 11228–11240 (2012).
46. D. Standage, M. Paré, Slot-like capacity and resource-like coding in a neural model of multiple-item working memory. *J. Neurophysiol.* **120**, 1945–1961 (2018).
47. R. Taylor, P. M. Bays, Efficient coding in visual working memory accounts for stimulus-specific variations in recall. *J. Neurosci.* **38**, 7132–7142 (2018).
48. X.-X. Wei, A. A. Stocker, A Bayesian observer model constrained by efficient coding can explain 'anti-Bayesian' percepts. *Nat. Neurosci.* **18**, 1509–1517 (2015).
49. D. Ganguli, E. P. Simoncelli, Efficient sensory encoding and Bayesian inference with heterogeneous neural populations. *Neural Comput.* **26**, 2103–2134 (2014).
50. K. Oberauer, H.-Y. Lin, An interference model of visual working memory. *Psychol. Rev.* **124**, 21–59 (2017).
51. S. Schneegans, P. M. Bays, Neural architecture for feature binding in visual working memory. *J. Neurosci.* **37**, 3913–3925 (2017).
52. F. Bouchacourt, T. J. Buschman, A flexible model of working memory. *Neuron* **103**, 147–160 (2019).
53. L. Matthey, P. M. Bays, P. Dayan, A probabilistic palimpsest model of visual short-term memory. *PLoS Comput. Biol.* **11**, e1004003 (2015).
54. P. M. Bays, Correspondence between population coding and psychophysical scaling models of working memory. [bioRxiv:10.1101/699884](https://doi.org/10.1101/699884) (11 July 2019).
55. H. Sompolinsky, H. Yoon, K. Kang, M. Shamir, Population coding in neuronal systems with correlated noise. *Phys. Rev.* **64**, 051904 (2001).
56. J. Jonides *et al.*, The mind and brain of short-term memory. *Annu. Rev. Psychol.* **59**, 193–224 (2008).
57. A. D. Baddeley, *Working Memory* (Oxford University Press, Oxford, United Kingdom, 1986).
58. P. Barrouillet, V. Camos, As time goes by: Temporal constraints in working memory. *Curr. Dir. Psychol. Sci.* **21**, 413–419 (2012).
59. K. Oberauer, S. Lewandowsky, S. Farrell, C. Jarrold, M. Greaves, Modeling working memory: An interference model of complex span. *Psychon. Bull. Rev.* **19**, 779–819 (2012).
60. S. Schneegans, R. Taylor, P. M. Bays, Stochastic sampling provides a unifying account of visual working memory limits. *Open Science Framework*. <https://osf.io/buxp9/>. Deposited 28 July 2020.