# Uncertainty Quantification in Working Memory via Moment Neural Networks

Hengyuan Ma[1], Wenlian Lu[1,2,3,4,5,6], Jianfeng Feng[1,2,3,4,7*]

1 Institute of Science and Technology for Brain-inspired Intelligence, Fudan University, Shanghai 200433, China
2 Key Laboratory of Computational Neuroscience and Brain-Inspired Intelligence (Fudan University), Ministry of Education, China
3 School of Mathematical Sciences, Fudan University, No. 220 Handan Road, Shanghai, 200433, Shanghai, China
4 Shanghai Center for Mathematical Sciences, No. 220 Handan Road, Shanghai, 200433, Shanghai, China
5 Shanghai Key Laboratory for Contemporary Applied Mathematics, No. 220 Handan Road, Shanghai, 200433, Shanghai, China
6 Key Laboratory of Mathematics for Nonlinear Science, No. 220 Handan Road, Shanghai, 200433, Shanghai, China
7 Department of Computer Science, University of Warwick, Coventry, CV4 7AL, UK
∗ jffeng@fudan.edu.cn

## ABSTRACT

Humans possess a finely tuned sense of uncertainty that helps anticipate potential errors, vital for adaptive behavior and survival. However, the underlying neural mechanisms remain unclear. This study applies moment neural networks (MNNs) to explore the neural mechanism of uncertainty quantification in working memory (WM). The MNN captures nonlinear coupling of the first two moments in spiking neural networks (SNNs), identifying firing covariance as a key indicator of uncertainty in encoded information. Trained on a WM task, the model demonstrates coding precision and uncertainty quantification comparable to human performance. Analysis reveals a link between the probabilistic and sampling-based coding for uncertainty representation. Transferring the MNN's weights to an SNN replicates these results. Furthermore, the study provides testable predictions demonstrating how noise and heterogeneity enhance WM performance, highlighting their beneficial role rather than being mere biological byproducts. These findings offer insights into how the brain effectively manages uncertainty with exceptional accuracy.

## 1 Introduction

Humans not only make decisions but also assess the degree of confidence associated with those decisions. This ability is crucial because it enables individuals to adapt their behavior based on the reliability of their judgments, enabling more effective navigation of uncertain environments. For instance, Accurately assessing confidence (or quantifying uncertainty) allows humans to prioritize cognitive resources, seek targeted information to resolve specific uncertainties, and effectively communicate their confidence levels to guide collaborative decision-making. Importantly, extensive studies have shown that humans possess the ability to generate a sense of uncertainty that accurately reflects the likelihood of errors across various tasks [1, 2, 3, 4, 5]. In other words, humans tend to feel less confident when they make larger mistakes, and more confident when their errors are smaller. This capability gives humans an advantage over artificial intelligence in tasks requiring trustworthiness and reliability, as deep networks often become overconfident in their predictions even when large mistakes occur [6]. However, the neural mechanisms underlying this precise uncertainty representation remain unclear. Elucidating the neural basis of uncertainty representation is a fundamental challenge in neuroscience [7]. Advancing this understanding also has implications for designing more trustworthy and interpretable artificial intelligence systems [8]. Notably, research in this domain has extensively utilized second-order moments for uncertainty quantification (UQ) [9, 10]. This raises an intriguing question: do humans similarly rely on second-order moments of neural activities for UQ?

There has been a long-standing debate on how neuron population activity represents uncertainty [11]. Two leading theories, probabilistic population coding [12] and sampling-based coding [13], offer different interpretations. Probabilistic population coding suggests that both the mean of the target variable and its uncertainty are captured by time-averaged neural responses. In contrast, sampling-based coding attributes the mean representation to the average response and uncertainty to response variability. While both frameworks successfully link neural activity to behavioral variability observed in experiments, neither fully captures the precise representation of uncertainty required to account for observed errors [3, 4, 5, 14].

Working memory (WM) is the basis of numerous high-level cognitive processes [15]. Understanding the neural mechanisms underlying uncertainty representation and quantification in WM is essential for explaining how humans manage uncertainty and generate confidence in various high-level tasks. Although WM is prone to errors due to both internal and external noise [16], humans are aware of the potential errors through a sense of uncertainty [14] and leverage this awareness to optimize their performance [3]. Additionally, significant fluctuations in WM error levels

have been observed [17, 1, 18], highlighting the importance of reliable UQ for recognizing deficiencies in WM content. Although probabilistic models of UQ have been proposed [19, 14], its neural mechanisms and a biologically plausible implementations, such as spiking neural networks (SNN), remain largely underexplored.

Ring attractor neural networks are prominent models for WM, supported by experimental evidence [20, 21, 22, 23]. These models apply a ring manifold to store a continuous feature, such as head direction [24]. However, they rely on highly structured synaptic connections [25, 26], which contradict the observed heterogeneity in neural tuning functions and synaptic connections [27, 22]. Additionally, WM errors are attributed to diffusion of the bump location, which fails to explain how humans quantify uncertainty in WM, as the variance of the decoded feature remains constant despite variations in actual memory error [1, 18].

Instead of manual design, there is a growing trend toward training recurrent neural networks for cognitive tasks [28, 29], which facilitates the discovery of diverse network configurations. However, many of these studies [30] rely on backpropagation through time for training [29], which is not biologically plausible. Recently, Darshan et al. (2022) proposed to train networks with synaptic heterogeneity using reservoir computing for WM [31]. While this approach is backpropagation-free and breaks the structured network connection, it introduces systematic drift, leading to significant coding errors without a mechanism to capture the associated uncertainty. Additionally, like many studies [30, 32], they use rate-based models, which fail to account for biologically realistic features, such as spike-based communication in neurons. Training more realistic models, like SNNs, remains much more challenging.

In this study, we address the above issues by employing moment neural networks (MNNs) [33]. Unlike rate-based neural models, which simplify SNNs by considering only the mean firing rate (mean), MNNs capture the complex nonlinear dynamics arising from both the mean firing rate and firing covariance (covariance) of leaky integrate-and-fire SNNs. Moreover, MNNs use differentiable moment activations (Eq. (4), *Methods*) instead of discrete spikes, allowing for easier training. MNNs have proven to be a valuable model for studying the neural basis of various cognitive functions. They have revealed how covariance influences coding precision and memory capacity in working memory (WM) [34], demonstrated how covariance encodes or extracts information during perceptual tasks [35], and shown how spiking neurons facilitate faster decision-making processes [36].

Our results show that, trained using the reservoir computing approach introduced in [31], MNNs effectively capture uncertainty in WM through covariance, achieving coding precision comparable to human performance in WM tasks, even with significant synaptic heterogeneity. We then propose a potential mechanism for uncertainty representation in WM, revealing a direct link between probabilistic population coding and sampling-based coding through the nonlinear coupling of mean and covariance. These findings, derived from the MNN, are validated on the SNN by transferring the trained MNN weights. This not only confirms that the MNN serves as a faithful substitute for the SNN but also achieves the first implementation of uncertainty quantification in WM via spike-based communications. Additionally, we present testable predictions on how factors such as noise, neural correlations, and heterogeneity influence WM performance, highlighting their beneficial roles during network training and inference. Overall, this study offers new insights into how humans manage uncertainty with remarkable precision through correlated neuronal fluctuations.

## 2  Working memory and uncertainty quantification

WM is one of the most widely studied cognitive functions, serving as the foundation for various cognitive processes, including learning and comprehension, attention and focus, goal-oriented behavior, and adaptability [37]. Previous studies have shown that WM is susceptible to errors, leading to variations in precision that are reflected in behavioral mistakes [16]. However, humans possess an awareness of these potential errors, which manifests as a sense of confidence about precision of remembered content [14]. Further research indicates that individuals can leverage this awareness to optimize their performance [3].

As shown in Fig. 1a, the WM task designed by Li et al. [14] involves participants first being presented with a cue that provides information about a feature (location) to be memorized. This is followed by a delay period during which the cue is removed. After the delay, participants are asked to report the remembered feature and indicate their uncertainty by adjusting the length of an arc centered at the reported location, representing a confidence interval for estimating the true position. To encourage participants to quantify uncertainty about the target location accurately, the experiment awarded points based on their responses. Points were granted only if the target location fell within the designated arc, with the score decreasing as the arc length increased. To maximize their scores, participants were expected to use shorter arcs when they were more confident. Broadly, responses can be categorized into four types, as illustrated in Fig. 1b. Ideal uncertainty quantification (UQ) occurs when the arc length is large for large errors and small for small errors. The other two cases represent overconfidence and underconfidence. The study reported a strong correlation (correlation coefficient of 0.6–0.7) between arc length and error (Fig. 5D in [14]), indicating that participants' reported uncertainty reliably reflected the error in their responses.
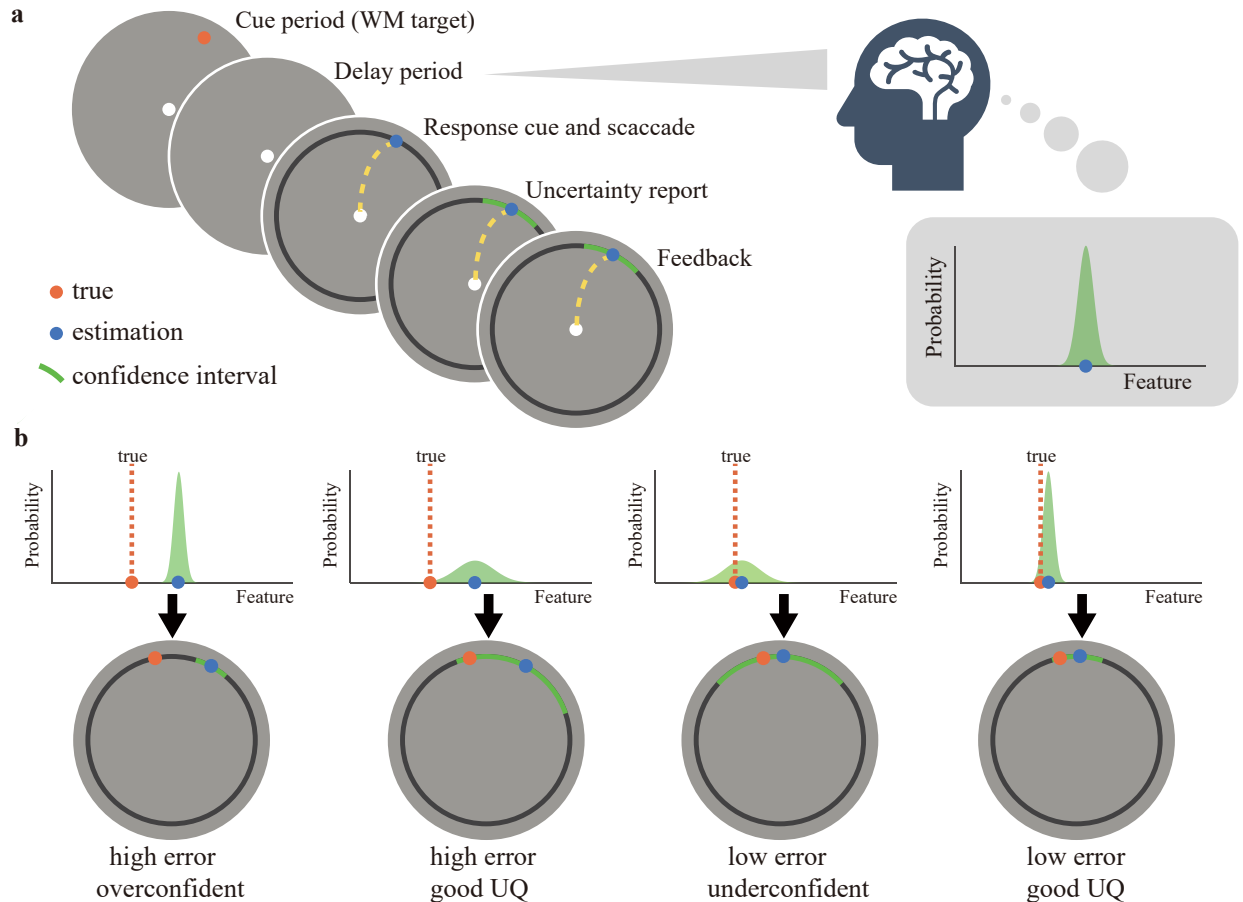
a

Cue period (WM target)

Delay period

Response cue and scaccade

Uncertainty report

Feedback

● true
● estimation
⌣ confidence interval

Probability

Feature

b

Probability — true — Feature
Probability — true — Feature
Probability — true — Feature
Probability — true — Feature

high error
overconfident

high error
good UQ

low error
underconfident

low error
good UQ

Figure 1: **Working memory task and its uncertainty quantification (UQ).** (a) In the task designed in [14], participants are required to remember the location indicated by the cue. After a delay period, they use saccades to indicate the remembered location on the ring and report their uncertainty with an arc. (b) Four representative cases of UQ results. Effective UQ should accurately reflect the magnitude of the error.

## 3 Neural models for working memory

### 3.1 Spiking neural network

Spiking neural networks (SNNs) (Eq. (3), *Methods*) are widely used to model the dynamics of biological neural systems underlying various cognitive tasks, as illustrated in the top of Fig. 2a. However, the discrete nature of spike trains presents challenges for model construction and analysis. To address this, spike trains can be simplified by extracting statistical features such as mean and covariance, as shown in Fig. 2b. When only the mean is considered, spiking neural networks can be represented as rate-based neural networks.

### 3.2 Rate-based neural network

Rate-based neural models, which serve as a simplification of SNNs, are widely used to model WM. Its dynamic is defined as

$$\tau \frac{\partial \boldsymbol{\mu}}{\partial t} = -\boldsymbol{\mu} + \phi(\bar{\boldsymbol{\mu}}), \tag{1}$$

where $\boldsymbol{\mu} \in \mathbb{R}^N$ represent the mean of the spike count per unit time (mean firing rate), $\phi(\cdot)$ is the element-wise nonlinearity, which is often set as Sigmoid or tanh function, $\tau$ is the membrane time constant, and $\bar{\boldsymbol{\mu}} \in \mathbb{R}^N$ corresponds to the mean of the input current for each neuron calculated as $\bar{\boldsymbol{\mu}} = W\boldsymbol{\mu} + \boldsymbol{\mu}_s$, where $W \in \mathbb{R}^{N \times N}$ is the synaptic connection matrix, and $\boldsymbol{\mu}_s$ is the external input current. We illustrate the rate-based neural network in the top of Fig 2b. To model WM, the weight connections in the rate-based neural model are set to be shift-symmetric, with long-distance
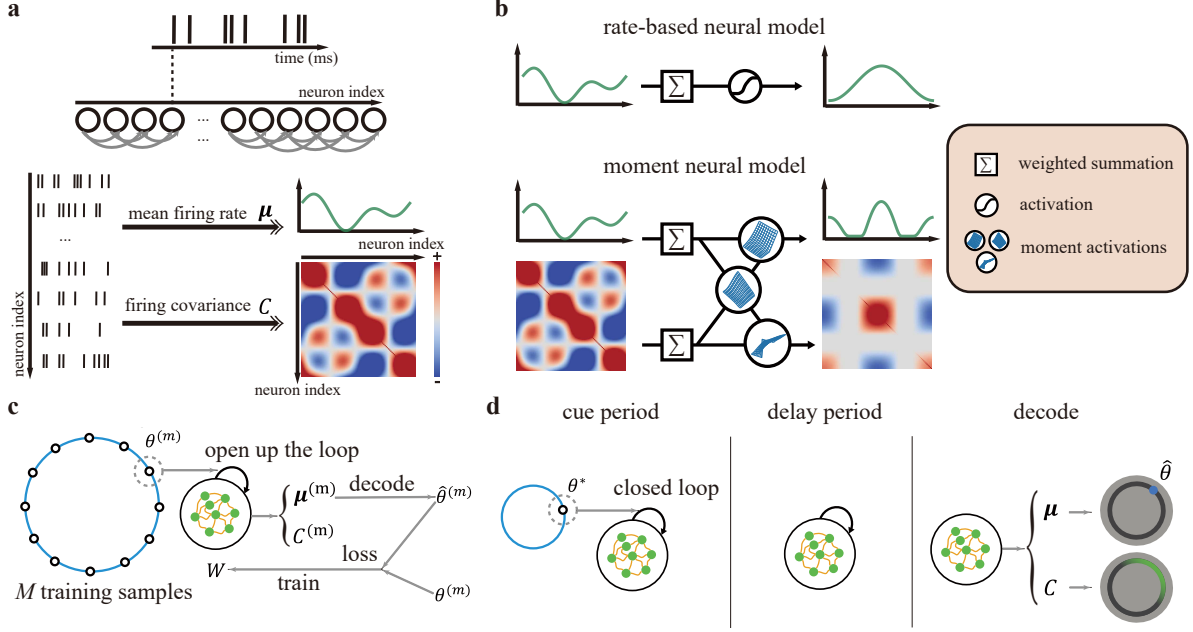
Figure 2: **Comparison of the spiking neural network (SNN), rate-based neural network, and moment neural network (MNN).** (a). (Top) An scheme of an SNN. (Bottom) Spike trains of a neuron population can be summarized by the mean firing rate of neurons and the firing covariance between each pair of neurons. (b). (Top) The rate-based neural model only considers the mean firing rate and its nonlinear evolution through an activation function. (Bottom) The MNN captures both the mean firing rate and firing covariance, with their nonlinear coupling during network evolution represented by the moment activations. (c) Training the synaptic connection of an MNN for the working memory task using reservoir computing. (d) The inference procedure is as follows: During the cue period, a feature encoded by external input is sent to the trained MNN. During the delay period, the external inputs are removed. After the delay period, the feature estimation and its uncertainty are decoded from the mean and covariance of the MNN, respectively.

inhibition and short-distance excitation. This configuration ensures that the system's fixed points are identical bumps that shift to form a ring manifold [20, 38]. The location of the bump encodes the feature being remembered. To incorporate internal neuronal noise, white noise is exerted to the network, causing diffusion of the bump location, with variance increasing linearly over time[39]. However, it cannot account for the variability in precision observed in WM experiments conducted under identical conditions [40, 17], hence inadequate for explaining the reliable UQ performance observed in humans.

We highlight a key distinction between spike-based biological neural systems and rate-based neural networks. In the former, each computational unit is governed by point-process (spikes), while in the latter, units are represented by scalar values. These scalars approximate the spike point process, representing its firing rate. However, this approximation overlooks the variability inherent in point processes and the nonlinear dynamics of such variability. Previous studies have demonstrated the significant role this variability plays in neural computation [41, 42].

### 3.3 Moment neural networks

A major challenge faced by rate-based neural models is difficulties in capturing realistic neuronal fluctuations (firing covariance) that are nonlinearly coupled with the the mean , which plays an important role in the neural dynamics [33, 43, 44]. To capture the nonlinear coupling of correlated fluctuations in the recurrent population of spiking neurons, we employ a model known as the moment neural network (MNN) [33] defined as

$$\begin{cases} \tau \frac{\partial \boldsymbol{\mu}}{\partial t} = -\boldsymbol{\mu} + \phi_\mu(\bar{\boldsymbol{\mu}}, \bar{C}) \\ \tau \frac{\partial C}{\partial t} = -C + \phi_C(\bar{\boldsymbol{\mu}}, \bar{C}) \end{cases} \tag{2}$$

where firing covariability $C \in \mathbb{R}^{N \times N}$ represent the covariance of the spike count per unit time, respectively. The covariance of the total synaptic current input $\bar{C} \in \mathbb{R}^{N \times N}$ is calculated as $\bar{C} = WCW^\top + \sigma_s^2 I$, where $\sigma_s$ is the noise level during the training phase. The moment activations $\phi_\mu$ and $\phi_C$ (defined in Eq. (4), *Methods*) describes the relationship between the input current statistics and the output spike train statistics in the leaky integrate-and fire (LIF)

4

spiking neural model. Unlike rate-based neural models which typically use heuristic activation functions such as tanh or sigmoid [45, 46, 47, 26], the nonlinearity of MNN are derived through a mathematical technique known as the diffusion approximation [48, 49] which faithfully captures the nonlinear coupling of mean and firing variability across populations of spiking neurons.

We highlight three key aspects of the MNN. First, unlike rate-based neural models (top of Fig. 2b), the MNN captures the nonlinear coupling between mean and covariance, as illustrated in the bottom of Fig. 2b. Second, the pattern of the covariance of the input synaptic current to the neural population $\bar{C}$ that emerges from the model is a result of the intrinsic dynamics of the recurrent circuit, rather than being determined by external input, as in rate-based neural networks, where the input is uncorrelated. Third, the nonlinear coupling between mean and covariance during MNN training allows the network to regulate both simultaneously in task performance. For instance, the error signal during training enables the model to adjust covariance to represent task-related uncertainty, a capability that networks based solely on mean cannot achieve.

## 4 Training the MNN for working memory

We trained an MNN to hold WM of a continuous periodic feature $\theta \in [0, 2\pi)$. During training, we selected $M$ training features $\theta^{(m)}, m = 1, \ldots, M$, each encoded as an external current $\boldsymbol{\mu}_s^{(m)}$. We trained the connections $W$ so that the mean of the fixed point encodes the value of $\theta^{(m)}$ under input $\boldsymbol{\mu}_s^{(m)}$. We applied a reservoir computing approach following [31], with $L_2$ regularization controlled by a factor $\alpha$ (Eq. (13) in *Methods*), which encourages weaker or sparser connections, making the network more energy-efficient. See *Methods* for details. The training procedure is illustrated in Fig. 2c.

Two key points should be emphasized. First, similar to the approach in [31], the connection matrix $W$ includes an untrainable random matrix component $J \in \mathbb{R}^{N \times N}$, which introduces unstructured heterogeneity, consistent with experimental findings that neuronal connection dynamics are highly heterogeneous [50, 51, 52]. The degree of this heterogeneity is controlled by a factor $g > 0$. While some studies have shown that structured heterogeneity can reduce error levels in WM [53], unstructured heterogeneity has been shown to often lead to a rapid decline in the network's computational capabilities [25]. This highlights the challenge of maintaining reasonable cognitive precision of our model. Second, the training loss (Eq. (13) in *Methods*) does not include the covariance $C$. Consequently, only the mean is directly supervised during training, while the covariance remains unsupervised. This distinguishes our MNN from many models in machine learning for UQ, which typically require supervision of both the mean and the (co)variance of the model's output through loss functions such as log-likelihood or evidence lower bound [54, 9, 55]. In contrast, the covariance in our MNN adapts to its inputs indirectly, through the nonlinear coupling between the mean and the covariance.

After training, we tested the network for the WM task on $L$ instances, as illustrated in Fig. 2d. We first input cue that convey the information of a target feature $\theta^*$ to the network during the cue period. Then, during the delay period, we removed the external current and introduced a new parameter $\sigma_{s,2}$ to control the noise level during the inference phase (see Eq. (14), *Methods*). After the delay period ends, we decoded the remembered feature $\hat{\theta}$ and its associated uncertainty level $\kappa$ from the network state as $(\boldsymbol{\mu}, C)$. We quantified the uncertainty level represented by the covariance based on four uncertainty metrics I-IV (Eq. (16) and Eq. (17), *Methods*), and we also record the estimation error $e$. An effective UQ is expected to show a strong positive correlation coefficient between the uncertainty $\kappa$ and the error $e$, which is estimated by Eq. (18), *Methods*.

## 5 Precise coding and reliable uncertainty quantification in the MNN

We analyzed the trained MNN as follows. Unlike previous studies where the bump attractor network's connection matrix is homogeneously designed [39, 20, 34], our MNN employs a connection matrix trained through reservoir computing, a biologically plausible algorithm [56], rather than being manually specified. This approach makes our model more biologically realistic: both the tuning curves (Fig. 3a) and the connections (Fig. 3b) are heterogeneous, aligning with experimental findings [52]. Interestingly, each neuron only responds to a narrow range of features, suggesting a sparse coding scheme often observed experimentally [57].

We then present the decoding results of the MNN at the inference phase. Six instances of the decoded results are shown in Fig. 3c, illustrating that error levels vary significantly across different instances. This variability is consistent with experimental observations of substantial fluctuations in the quality of WM representations within an individual [40, 17]. More importantly, we observe that larger errors are associated with higher covariance levels, suggesting that covariance effectively captures uncertainty. Furthermore, we divided all $L$ test instances into three groups based on the uncertainty
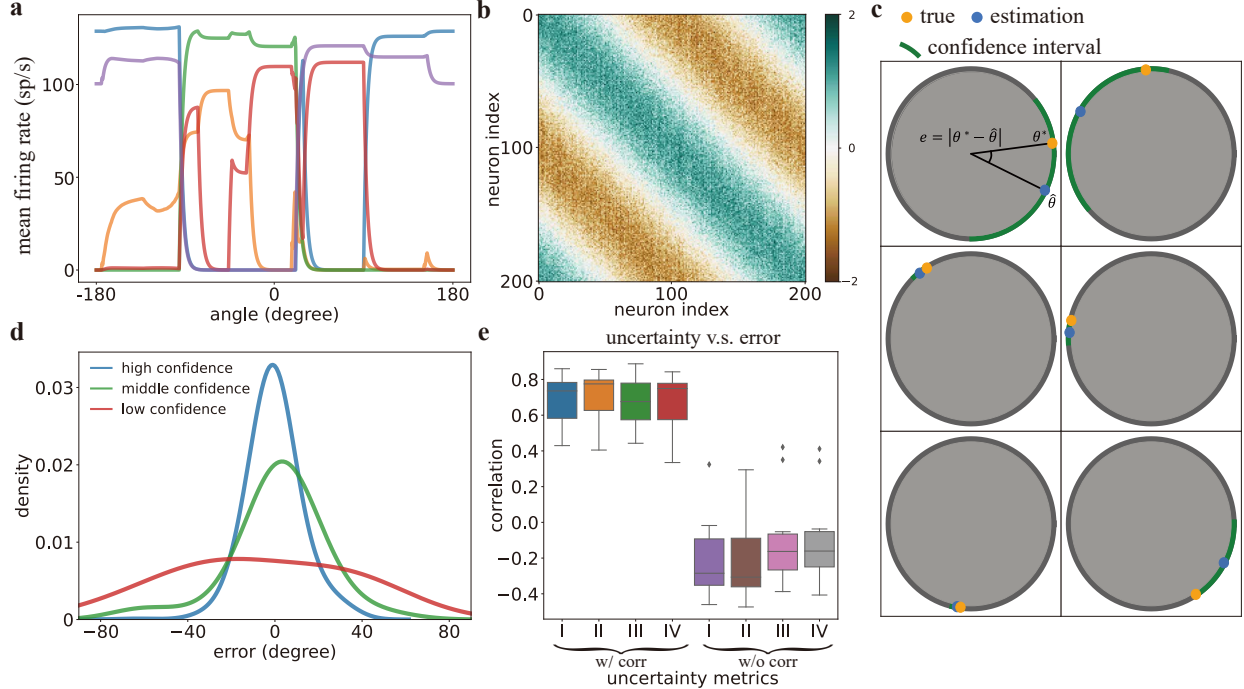
Figure 3: **The performance of moment neural network (MNN) trained on the working memory task.** (a) The tuning curves of five neurons. (b) The trained weights of the network. (c) Six instances of the true input variables and the decoded feature and confidence interval decoded from the network. The arc length of the confidence interval is proportional to the square root of the first eigenvalue of the decoded covariance matrix $\hat{C}_z$, see **Network inference** section in *Methods*. (d) The error distribution across three groups of instances, divided based on the level of uncertainty (using the uncertainty metrics I defined Eq. (16), *Methods*): top 25% uncertainty (low confidence), top 25-50% uncertainty (middle confidence), and the remaining instances (high confidence). Corresponding results of uncertainty metrics II-IV are shown in the *Supplementary Information*. (e) The correlation between uncertainty (calculated using four indicators, I-IV) and the error, calculated under two conditions: one where the correlation between neuron activities is maintained (w/ corr) and one where the correlation is clamped to zero (w/o corr).

level: the low-confidence group corresponds to the top 25% of uncertainty indicator $I$, the middle-confidence group corresponds to the 25%-50% range of uncertainty indicator $I$, and the remaining instances correspond to the high-confidence group. As shown in Fig.3d, the low-confidence group exhibits the heaviest-tailed error distribution, followed by the middle-confidence group. This observation aligns with the experimental results in [18, 40]. Additionally, the precision of our model also reaches comparable level reported in [18, 40, 14].

We then systematically evaluated the UQ capability of the MNN. As shown in Fig. 3e, the uncertainty levels $\kappa$ calculated from the I-IV indicators are all strongly positively correlated with the error $e$, with correlation coefficients comparable to those reported in [14]. This suggests that the covariance in the MNN effectively captures uncertainty.

To examine the role of neuron firing correlation in UQ, we clamped all pairwise correlations between neurons in the covariance matrix to zero in the MNN and repeated the experiments as an ablation study. As shown in Fig. 3e, we observed a significant decrease in the correlation coefficient $\rho$ between the error $e$ and each of the four indicators. This finding suggests that the correlation between neural activities is crucial for the neural system's ability to quantify uncertainty. It aligns with studies showing through animal experiments that correlated fluctuations between neurons are essential for population coding [58].

## 6    Mechanism of the uncertainty quantification

We analyzed the patterns of the fixed points, as shown in Fig. 4a, where the mean and covariance of three fixed points are presented. Each pattern has been centered for comparison, and the five patterns are arranged in increasing order of covariance level. The mean shows an imperfect bump, while the covariance displays a square-like pattern. More importantly, the bump width (defined as the number of neurons with firing rates exceeding 5% of the peak) increases
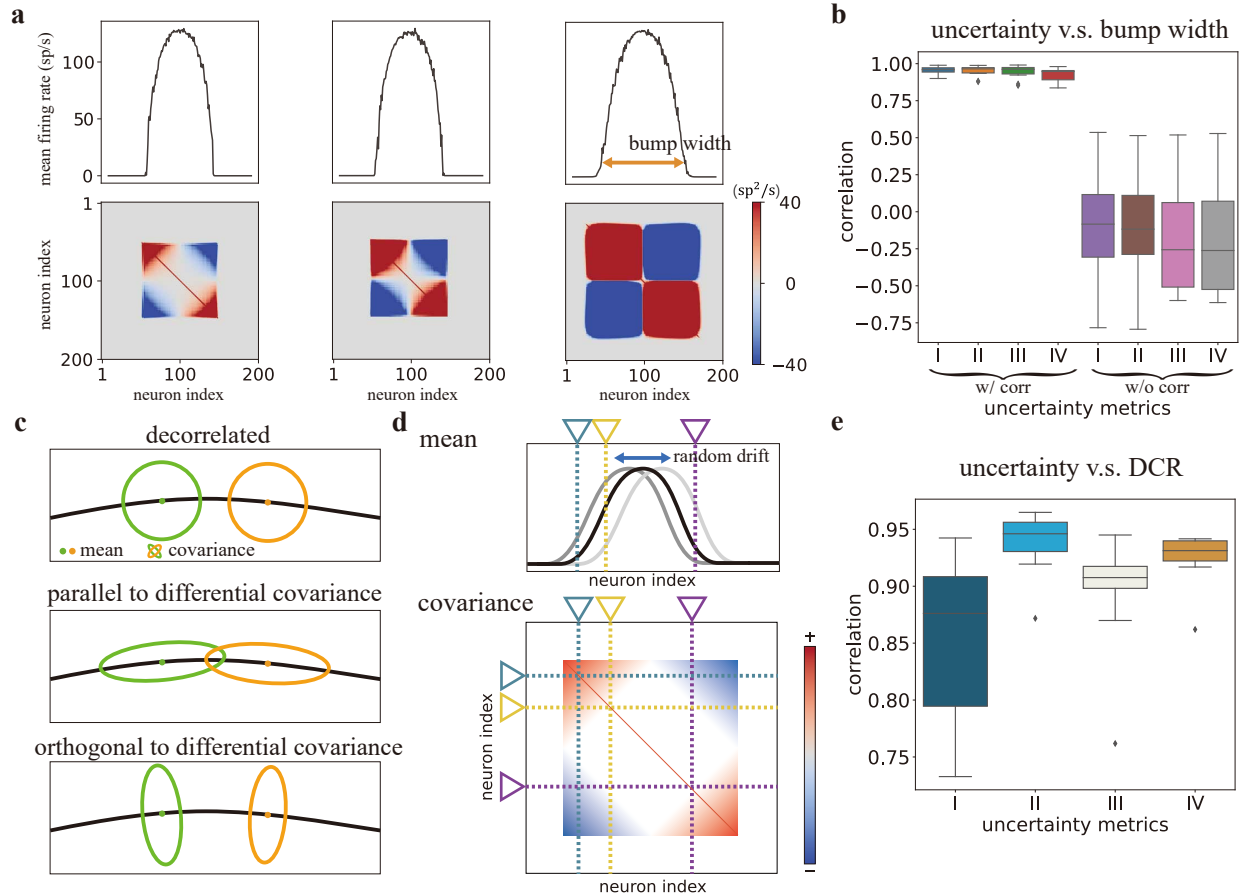
Figure 4: **Mechanism of uncertainty quantification in the moment neural network (MNN) for working memory.**
(a) Several typical fixed-point patterns of mean firing rate (top) and firing covariance (bottom) produced by the MNN, with the bump width increasing from left to right. (b) The correlation between the four uncertainty indicators (I-IV, Eq. (16)-Eq. (17), *Methods*) and the bump width, calculated under two conditions: one where the correlation between neuron activities is maintained (w/ corr), and one where the correlation is clamped to zero (w/o corr). (c) A schematic showing how differential covariance affects the neural coding of variables with different values. (d) The random drift of the bump generates differential covariance components in the firing covariance. (e) The correlation between uncertainty and the differential covariance ratio (DCR).

with the covariance level. To systematically verify this, we calculated the correlation coefficient between bump width and the four uncertainty indicators. As shown in Fig. 4b, bump width is positively correlated with all four uncertainty indicators. In contrast, when firing correlations are clamped to zero, the correlation drops dramatically. This suggests that the uncertainty represented by the covariance is related to the bump width, which has been considered a form of probabilistic population coding for uncertainty [24]. In the probabilistic population coding framework, the mean encodes both the target variable and its associated uncertainty [12, 59]. For instance, the bump amplitude or width [24] represent uncertainty, where a smaller bump amplitude or larger bump width reflects greater uncertainty. In contrast, the sampling-based coding theory suggests that variability in neural activity encodes uncertainty, with higher variability corresponding to higher uncertainty. Our findings suggest that two theories capture distinct aspects of neural activity, and are linked through the nonlinear interaction between mean and covariance.

Next, we investigated how the nonlinear coupling between mean and covariance help to UQ. We propose that covariance reliably captures uncertainty through the following mechanism: when coding accuracy for an instance is low, the differential covariance component of the covariance increases, which has the effect of amplifying the level of uncertainty. Differential covariance is the covariance component that changes the neural responses across trials along the direction parallel to the tuning curve derivative hence shrinks its information it code [60, 61, 34]. As illustrated as Fig. 4c, the differential covariance makes two codes of feature less distinguishable, reducing the coding accuracy, which has also been empirically verified in [34]. If the coding accuracy of an instance is low, the bump location tends to drift from
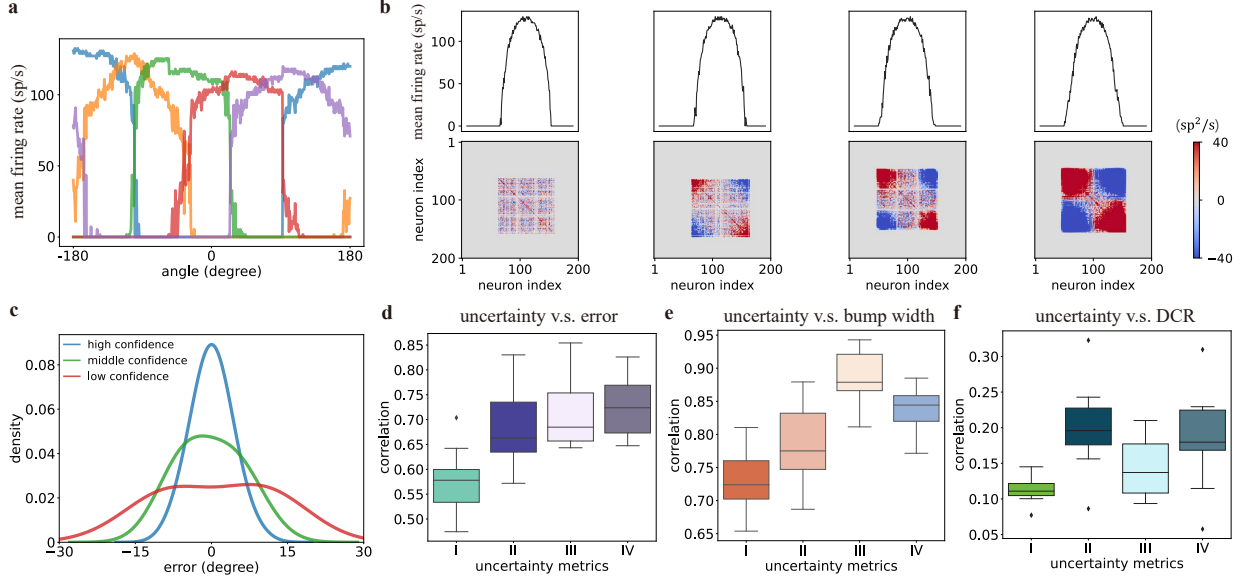
7

Figure 5: **The performance and spiking neural network (SNN) using the weights of the trained the moment neural network (MNN).** (a) The tuning curves of five neurons. The trained weights of the network. (b) Several typical fixed-point patterns of mean firing rate (top) and firing covariance (bottom) produced by the SNN, with the bump width increasing from left to right. (c) The error distribution across three groups of instances, divided based on the level of uncertainty (using the uncertainty metrics I defined Eq. (16), *Methods*): top 25% uncertainty (low confidence), top 25-50% uncertainty (middle confidence), and the remaining instances (high confidence). Corresponding results of uncertainty metrics II-IV are shown in the *Supplementary Information*. (d) The correlation between uncertainty (calculated using four indicators, I-IV) and the error. (e) The correlation between the four uncertainty indicators (I-IV, Eq. (16)-Eq. (17), *Methods*) and the bump width. (f) The correlation between uncertainty and the differential covariance ratio (DCR).

its original position, induces correlated neuronal fluctuation, as shown in the top of Fig. 4d. We hypothesize that the covariance state accounts for the effects of these fluctuations, leading to the emergence of a differential covariance structure, where neurons on the same side of the bump are positively correlated, and neurons on opposite sides show negative correlation, as illustrated at the bottom of Fig. 4d. Consequently, lower coding accuracy is associated with higher differential covariance components within the covariance. To verify this hypothesis, we introduce an indicator called the differential covariance ratio (DCR), which quantifies the amount of differential covariance (*Methods*). We calculated the correlation between DCR and the four uncertainty indicators. As shown in Fig. 4e, DCR is highly positively correlated with all four uncertainty indicators, suggesting that differential covariance is the primary source of uncertainty, hence supporting our hypothesis.

Building on the intuitive interpretation above, we provide a rigorous proof using a simplified model to demonstrate how a neural system trained with a loss function that only supervises the mean can also learn the true output variance, emphasizing how the mean-covariance coupling contributes to effective uncertainty quantification. See Thm. S1 and Thm. S2 in the *Supplementary Information*.

# 7 Verification on spiking neural network

We validated our findings on the SNN by transferring weights from the trained MNN to an SNN applying leaky integrate-and-fire neural model (Eq. (3), *Methods*). As shown in Fig. 5a, the tuning curves of the SNN were highly heterogeneous, consistent with those of the MNN. Again, each neuron responds to a narrow range of features, indicating a sparse coding scheme. The steady-state mean firing rates and firing covariances of the SNN (Fig. 5b) also exhibited qualitatively similar patterns to those observed in the MNN (Fig. 4a). An analysis of the SNN's decoding results revealed that the error distribution (Fig. 5c) and UQ performance (Fig. 5d) were comparable to those of the MNN. These results confirm that the MNN effectively captures the mean and covariance dynamics of the SNN. Therefore, the MNN is a reliable and efficient alternative to directly training SNN for cognitive tasks.

Furthermore, we observed a strong correlation between bump width and uncertainty in the SNN (Fig. 5e), indicating the presence of mean-covariance coupling. Additionally, we calculated the correlation between uncertainty and the DCR in Fig. 5f, which remained moderately positive, though lower than in the MNN (Fig. 4e). This reduction may be attributed to fluctuations in the estimated covariance of the SNN, which likely diminish certain components of the differential covariance. These findings suggest that the proposed UQ mechanism and mean-covariance coupling identified in the MNN remain valid in the SNN.

We note several advantages of the MNN over the SNN: First, different from the indifferentiable spike trains in the SNN, the moment activations of the MNN (Eq. (4), *Methods*) are differentiable, enabling easier training. Second, the MNN requires only a single pass to estimate the mean and covariance of the network state without fluctuations. In contrast, the SNN must be run multiple times to suppress fluctuations and achieve accurate estimates, leading to higher computational costs. Third, the MNN allows investigations into the role of covariance by clamping the off-diagonal elements of the covariance matrix to zero. Performing such manipulations in the SNN is considerably more challenging. These advantages establish the MNN as a more flexible framework for investigating neural mechanisms underlying cognitive functions.

## 8 Further analysis

We examined how factors including noise level and heterogeneity influence the accuracy and UQ results of the MNN. This analysis not only enhances the plausibility of our model, but also provides a deeper understanding of how these factors shape the neural processes underlying WM, particularly in terms of whether noise and heterogeneity—two inevitable properties of biological neural systems—are detrimental or beneficial to WM performance. The following results applied the uncertainty metrics I (Eq. (16), *Methods*), while the results of uncertainty metrics II-IV are shown in the *Supplementary Information*.

As shown in Fig. 6a, as the number of training samples $M$ increases, the quality of UQ initially improves significantly, then levels off. Previous work has shown that $M = 12$ is sufficient for training a bump attractor network [31], with fewer samples proving inadequate. This suggests that the network requires a sufficient number of samples to establish the underlying low-dimensional coding space of the feature $\theta$, with further increases in $M$ providing marginal improvement. Additionally, the correlation between uncertainty and bump width is less affected by $M$, suggesting that the coupling between the mean and covariance is less influenced by network training and is instead an intrinsic property of the MNN. Moreover, the error level continuously decreases as $M$ increases (Fig. 6c), which is consistent with the intuition that more training samples lead to better performance.

We then investigated the effect of the population size $N$. As shown in Fig. 6d, $N$ has a similar effect on UQ performance as the number of training samples $M$, i.e., the network requires a sufficient number of neurons (about 200) for stable UQ performance, with further increases in $N$ providing only marginal improvement. In contrast to $M$, increasing $N$ enhances the correlation between uncertainty and bump width (Fig. 6e). This suggests that enlarging the neuron population strengthens the coupling between mean and covariance. Additionally, the error level decreases as $N$ increases (Fig. 6f), likely because total information increases sub-linearly with population size, consistent with previous findings [62].

We analyzed the effect of the regularization factor $\alpha$ (Eq. (13), *Methods*). As shown in Figs. 6g and 6i, there is an optimal value of $\alpha$ that yields the best UQ performance and minimal error level, respectively. This suggests that appropriate regularization is necessary for optimal network performance. Additionally, as shown in Fig. 6h, the regularization parameter $\alpha$ has less impact on the correlation between uncertainty and bump width, further supporting the idea that the coupling between mean and covariance is an intrinsic property of the MNN, largely unaffected by training.

Next, we analyzed the effect of noise. As shown in Figs. 7a and 7c, the noise level during training, $\sigma_s$, significantly improves both UQ performance and accuracy. This suggests that noise plays a key role in error-awareness, underscoring its benefits for learning in the brain and highlighting its critical role in training neural networks for cognitive functions. The likely mechanism is that noise enhances the robustness of the network. Additionally, as shown in Fig. 7b, $\sigma_s$ also increases the correlation between uncertainty and bump width, thereby strengthening the coupling between mean and covariance. As shown in Figs. 7d-e, the noise level during inference, $\sigma_{s,2}$, increases the correlation between uncertainty and both error and bump width. Furthermore, increasing $\sigma_{s,2}$ raises the error level, as expected (Fig. 7f). These results show that the MNN can maintain strong UQ performance across different levels of external noise, even when the training noise level is fixed. This suggests that the MNN generalizes its UQ ability to various noise conditions, a property that may be crucial for adaptation to new environments based on learned knowledge, as seen in humans and other species. We also provide a theoretic interpretation of such generalization ability in Thm. S1 in the *Supplementary Information*.
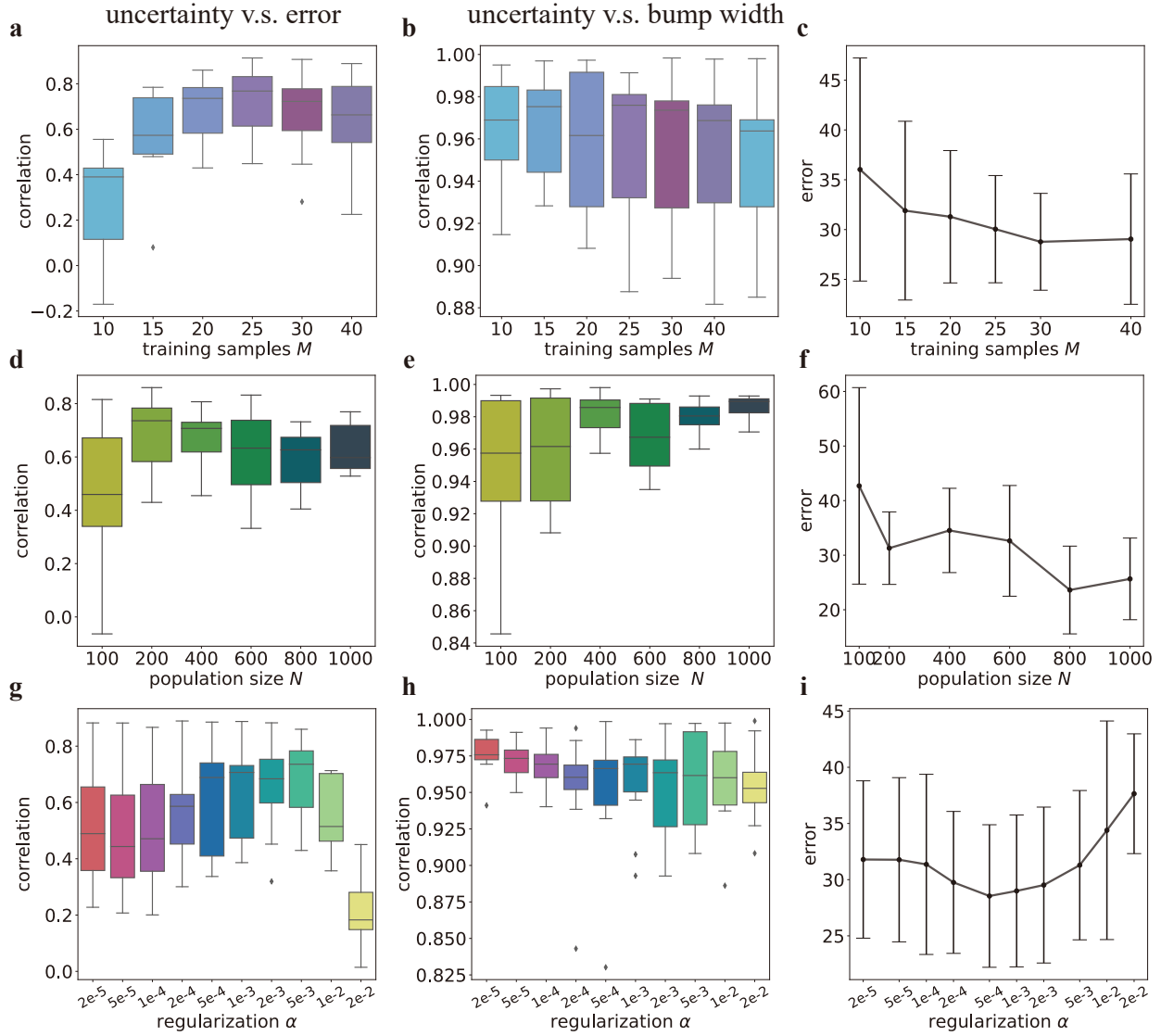
Figure 6: **Effect of training samples $M$, population size $N$, and regularization $\alpha$ on uncertainty quantification, mean-covariance coupling, and precision in working memory (WM) tasks.** (a-c) We conduct the same experiments as in Fig. 3 with varying training samples $M$, while keeping other parameters constant. (d-f) The same experiments as in (a-c), but with different population sizes $N$. (g-i) The same experiments as in (a-c), but with different regularization values $\alpha$.
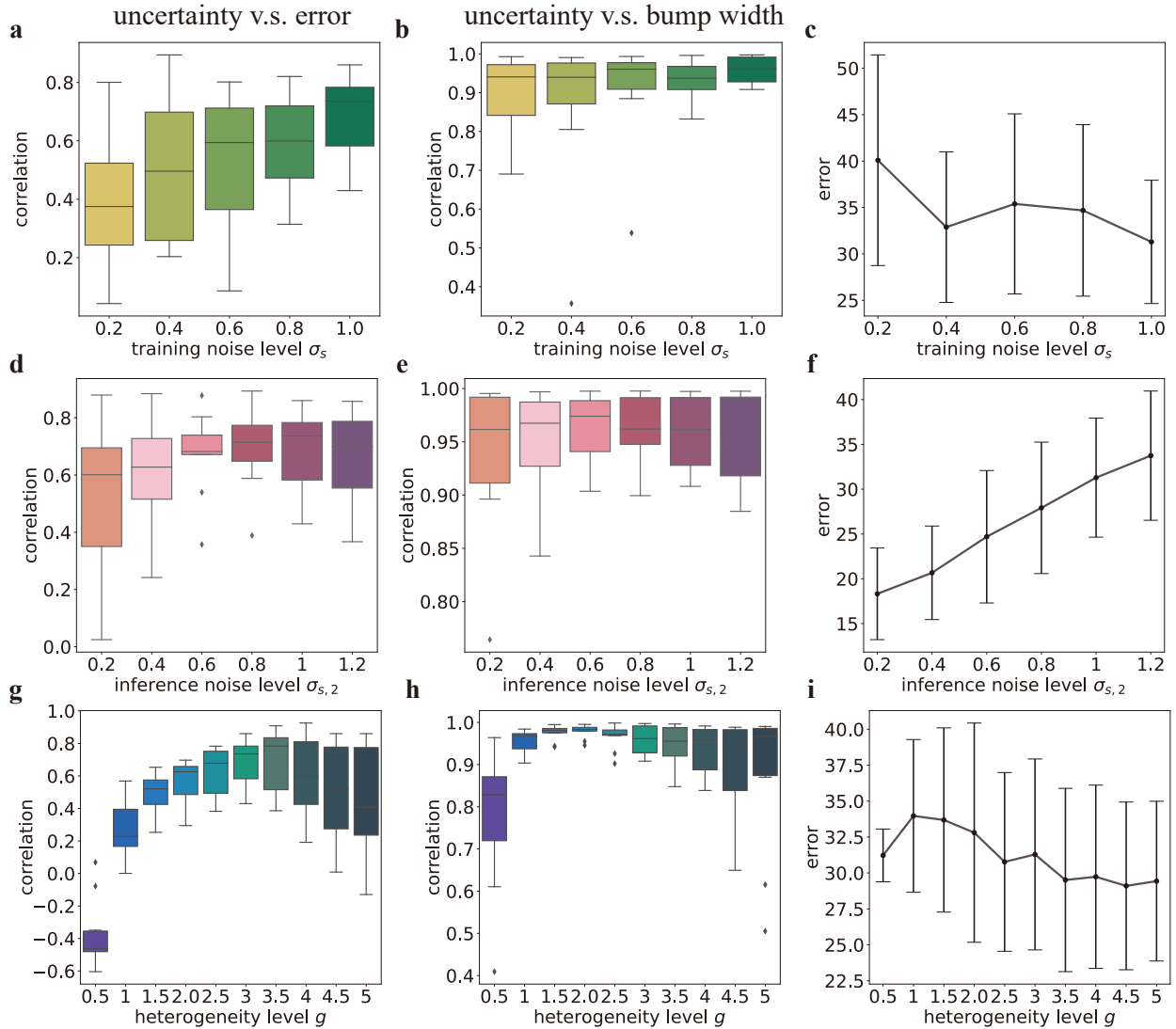
Figure 7: **Effect of level of noise at training phase $\sigma_s$, level of noise at inference phase $\sigma_{s,2}$, and heterogeneity $g$ on uncertainty quantification, mean-covariance coupling, and precision in working memory (WM) tasks.** (a-c) We conduct the same experiments as in Fig. 3 with varying level of noise at training phase $\sigma_s$, while keeping other parameters constant. (d-f) The same experiments as in (a-c), but with varying level of noise at inference phase $\sigma_{s,2}$. (g-i) The same experiments as in (a-c), but with varying level of heterogeneity $g$.

11

Finally, we investigated the effect of heterogeneity. As shown in Figs. 7g-h, when heterogeneity $g$ is low, both the correlation between uncertainty and error, and the correlation between uncertainty and bump width are also low, resulting in poor UQ. This suggests that some degree of heterogeneity is necessary for the mean-covariance coupling. Interestingly, as shown in Fig. 7i, the error level first increases and then decreases as heterogeneity increases. This result is similar to that observed in [53], where heterogeneity was introduced into a network to improve coding precision. The difference is that in our study, heterogeneity is unstructured, whereas in [53] it is spatially structured. Thus, our results extend this earlier work by showing that the model can benefit from unstructured heterogeneity, which is more commonly observed in biological neural systems.

## 9  Discussion

Humans are prone to making mistakes in their working memory (WM), yet they can generate a sense of confidence that accurately reflects the level of error. In this study, we investigate the neural mechanisms underlying uncertainty quantification (UQ) using the moment neural network (MNN), which links probabilistic population coding and sampling-based coding theories through its nonlinear coupling of the mean and covariance. Trained with reservoir computing to encode a periodic feature, the MNN demonstrates UQ performance comparable to that of humans (Fig. 3). These results are verified in an SNN using the trained MNN weights, resulting in a spike-based, heterogeneous, and sparse implementation of WM with its uncertainty quantification. Furthermore, we show that noise and heterogeneity, often seen as detrimental, actually enhance WM performance within a certain range (Figs. 7d-i), suggesting they are essential for task performance rather than mere byproducts.

We establish a direct connection between probabilistic population coding and sampling-based coding theories, two prominent theories of how neuronal activity encodes uncertainty. As illustrated in Fig. 4, we demonstrate that the bump width and the level of uncertainty calculated from the covariance are tightly coupled. The bump width serves as an indicator of uncertainty within the framework of probabilistic population coding, while the covariance acts as the uncertainty indicator in sampling-based coding. This coupling also explains why both theories are able to account for a wide range of experiments. More importantly, we propose a hypothesis for how covariance faithfully captures uncertainty: the error arises from the random drift of the bump, which amplifies the differential covariance component within the covariance, thereby increasing the level of uncertainty. Supporting this hypothesis, we find that the level of covariance is strongly correlated with the amount of differential covariance (Fig. 4e).

This study generates several testable predictions. First, the correlation between neural firing is crucial for both effective UQ (Fig. 3e) and mean-covariance coupling (Fig. 4b). This can be tested by decorrelating neuronal recording data to assess whether UQ performance is impaired and whether mean-covariance coupling disappears. Second, increasing the size of the neural population enhances UQ performance (Fig. 6d), although this effect may saturate once the population size exceeds a certain threshold. This prediction could be tested by comparing the number of neurons across different species and evaluating their respective quantification abilities. Third, the system may benefit from learning under noisy conditions (Fig. 7a-c). This prediction could be tested by designing experiments where subjects learn a task under varying levels of noise, then assess how performance changes with different levels of training noise. Fourth, network heterogeneity may enhance both UQ and coding precision (Figs. 7g-i). This hypothesis could be tested by quantifying heterogeneity across subjects and correlating it with performance.

The trade-off between model realism and feasibility for study and analysis is a significant challenge in neuroscience. Rate-based neural networks are widely used due to their ease of analysis and training; however, they sacrifice many biologically realistic features, such as spike communication and the intrinsic nonlinear dynamics of firing covariance, which play a crucial role in uncertainty quantification, as demonstrated in this work. On the other hand, spiking neural networks (SNNs) are more biologically accurate but much harder to analyze and train. The MNN applied in this work strikes a balance, combining the benefits of both approaches: the moment activations are differentiable, making the MNN easy to analyze and train like a rate-based neural network, while it captures the first two moments of the SNN, thus serving as a faithful substitute. Similar to our MNN, Echeveste et al. (2020) analyzed the dynamics of the first two moments of a neural network [32]. However, their work focuses on a rate-based neural network with additive noise, meaning that their network weights cannot be directly transferred to an SNN to replicate similar performance. In contrast, our verification experiments (Fig. 5) demonstrate that by training an MNN via reservoir computing and transferring its weights to an SNN, we can achieve performance in WM tasks that is comparable to the MNN.

Our results shed light on how the brain learns for UQ. Notably, the MNN achieves high UQ performance despite its loss function (Eq. (13), *Methods*) does not explicitly incorporate firing covariance, the indicator of the uncertainty. This means that only the mean is directly supervised for the task, while the covariance is indirectly adjusted through the nonlinear coupling between the mean and covariance. This provides an example of how higher-order statistics of a system can be regulated through first-order statistics. These findings also offer valuable insights for designing UQ

algorithms in machine learning, where UQ plays a critical role in applications such as active learning, reinforcement learning, domain adaptation, and security [6]. Specifically, in a supervised learning setting, it suggests an approach where only the mean (prediction) of a model output is supervised, while the output covariance adapts unsupervised, allowing the covariance to learn to quantify the uncertainty associated with the mean.

Aitchison et al. (2021) suggest that uncertainty can be encoded in synaptic weights and regulate the synaptic adaptation rate [63]. This implies that the uncertainty encoded in the MNN may regulate neuronal adaptation, playing a role in broader learning processes. Additionally, an intriguing direction for future research would be to explore how uncertainty is jointly or cooperatively encoded in the synaptic weights and neuron states and how they interact.

More broadly, the MNNs can be applied to model more complex cognitive processes. Since WM forms the foundation for higher-level functions, errors in WM can propagate to these processes, influencing the observable performance. Our study could be extended by analyzing how the firing covariance and the associated uncertainty in WM propagates to higher cognitive functions, offering insights into how humans quantify uncertainty in those tasks. In conclusion, our study provides opportunities for a deeper understanding of UQ and offers valuable insights into how biological systems manage uncertainty in noisy environments.

## Methods

**Spiking neural networks (SNNs).** We applied leaky integrate-and-fire (LIF) neural model as for the SNN in the study. The dynamics of an SNN with $N$ LIF neurons is defined as

$$
\begin{cases}
\frac{d\mathbf{u}}{dt} &= -L\mathbf{u} + \mathbf{I}, \\
\tau_s \frac{d\mathbf{I}}{dt} &= -\mathbf{I} + W\mathbf{s},
\end{cases}
\tag{3}
$$

where $\mathbf{u} \in \mathbb{R}^N$ is the membrane potential each neuron, $L \in \mathbb{R}$ is the leak conductance, $W \in \mathbb{R}^{N \times N}$ is the synaptic weights, $\mathbf{I} \in \mathbb{R}^{\mathbf{N}}$ is the synaptic currents, $\tau_s$ is the synaptic time constant, and $\mathbf{s} \in \mathbb{R}^N$ is the spike train generated by the neuron population. For $i = 1, \ldots, N$, when $u_i$ reaches the firing threshold $u_{\text{th}}$, the neuron emits a spike ($s_i = 1$) that is transmitted to connected neurons. Following the spike, $u_i$ is reset to the resting potential $u_{\text{res}}$ and enters a refractory period of duration $T_{\text{ref}}$. If $u_i$ does not reach $u_{\text{th}}$, no spike is emitted ($s_i = 0$). Throughout this work, we set neuron parameters to be $u_{\text{th}} = 20$ mV, $u_{\text{res}} = 0$ mV, $T_{\text{ref}} = 5$ ms, $L = 0.05$ ms$^{-1}$, and $\tau_s = 10$ ms.

**Moment neural networks (MNNs).** We employed a neural model called as the MNN, which is derived from the LIF neural model [64, 33]. Different from common rate-based neural model using elementiwise nonlinearity such as tanh and sigmoid. The moment activations $\phi_\mu \in \mathbb{R}^N \times \mathbb{R}^{N \times N} \to \mathbb{R}^N$ and $\phi_C \in \mathbb{R}^N \times \mathbb{R}^{N \times N} \to \mathbb{R}^{N \times N}$ together map the mean $\bar{\boldsymbol{\mu}}$ and covariance $\bar{C}$ of the steady-state input current $\mathbf{I} = W\mathbf{s}$ to that of the output spikes according to [64, 33]

$$
\phi_\mu(\bar{\boldsymbol{\mu}}, \bar{C})_i = \left(T_{\text{ref}} + \frac{2}{L} \int_{I_{\text{ub},i}}^{I_{\text{lb},i}} g(x) dx\right)^{-1},
$$

$$
\phi_C(\bar{\boldsymbol{\mu}}, \bar{C})_{ij} = \begin{cases}
\frac{8}{L^2} \phi_\mu(\bar{\boldsymbol{\mu}}, \bar{C})_i^3 \int_{I_{\text{ub},i}}^{I_{\text{lb},i}} h(x) dx, & i = j \\
\left(\frac{\partial \phi_\mu(\bar{\boldsymbol{\mu}}, \bar{C})}{\partial \bar{\boldsymbol{\mu}}}\right)_{ii} \left(\frac{\partial \phi_\mu(\bar{\boldsymbol{\mu}}, \bar{C})}{\partial \bar{\boldsymbol{\mu}}}\right)_{jj} \bar{C}_{ij}, & i \neq j
\end{cases},
\tag{4}
$$

where $g(x)$ and $h(x)$ are Dawson-like functions and defined as

$$
g(x) = e^{x^2} \int_{-\infty}^{x} e^{-u^2} du, \quad h(x) = e^{x^2} \int_{-\infty}^{x} e^{-u^2} [g(u)]^2 du.
\tag{5}
$$

and integration bounds are calculated as

$$
I_{\text{ub},i} = \frac{u_{\text{th}} L - \bar{\mu}_i}{\sqrt{L C_{ii}}}, \quad I_{\text{lb},i} = \frac{u_{\text{res}} L - \bar{\mu}_i}{\sqrt{L C_{ii}}}
\tag{6}
$$

**Simulation of MNNs.** All simulations were performed using the Euler method with a timestep $dt = 0.1\tau$ and $\tau = 1$ ms. The simulations were run on a single NVIDIA 3090 GPU. For the results presented in Fig. 3-7, each test was repeated 10 times. An efficient numerical algorithm was used to implement the moment neural network, as described in [65].

**Network connection and readout matrix.** The connection weight matrix is supposed to have a low-rank structure with additive quenched noise that brings heterogeneity:

$$
W = W_{\text{fb}} W_{\text{out}}^\top + gJ,
\tag{7}
$$

13

where the random matrix $J_{ij} \sim \mathcal{N}(0, 1/N)$ is fixed after being generated, and $g > 0$ is the heterogeneity level. The matrix $W_{\text{fb}} \in \mathbb{R}^{2 \times N}$ is fixed as $W_{1j,fb} = \cos(2\pi j/N), W_{2j,fb} = \sin(2\pi j/N)$, following [31], and the only trainable parameters are $W_{\text{out}} \in \mathbb{R}^{2 \times N}$, which is also the readout matrix. We set the heterogeneity $g = 3$ throughout the work, except in Fig. 7, where we vary $g$ from 0.5 to 5 to investigate its effect on the performance of the MNN. We set $N = 200$ throughout the work, except in Fig. 6, where we vary the population size $N$ from 100 to 1000 to investigate its effect on the performance of the MNN.

**Encoding and decoding in the MNN.** To input a cue of feature $\theta$ to the MNN, we encode it as the external current as $\boldsymbol{\mu}_s = W_{\text{fb}}\mathbf{z}$ with

$$\mathbf{z}_1 = A\cos(\theta), \quad \mathbf{z}_2 = A\sin(\theta), \tag{8}$$

where we set $A = 1.2$ throughout the work. When the state of the network as $(\boldsymbol{\mu}, C)$, the decoded mean and covariance are

$$\hat{\boldsymbol{\mu}}_z = W_{\text{out}}^\top \boldsymbol{\mu}, \quad \hat{C}_z = W_{\text{out}}^\top C W_{\text{out}}. \tag{9}$$

The feature is decoded by finding the angle $\hat{\theta}$ satisfying

$$\cos(\hat{\theta}) = \frac{\hat{\boldsymbol{\mu}}_{z,1}}{\sqrt{\hat{\boldsymbol{\mu}}_{z,1}^2 + \hat{\boldsymbol{\mu}}_{z,2}^2}}, \quad \sin(\hat{\theta}) = \frac{\hat{\boldsymbol{\mu}}_{z,2}}{\sqrt{\hat{\boldsymbol{\mu}}_{z,1}^2 + \hat{\boldsymbol{\mu}}_{z,2}^2}}. \tag{10}$$

To quantify the uncertainty associated with the decoded covariance $\hat{C}_z$ by considering the entropy of the Gaussian distribution $\mathcal{N}(\hat{\boldsymbol{\mu}}_z, \hat{C}_z)$, which is a common indicators of uncertainty of a distribution. Please see the **uncertainty metrics** section in *Methods*.

**Network training.** When uniformly select $M$ features from $[0, 2\pi)$ denoted as $\theta^{(m)}, m = 1, \ldots, M$. For each feature $\theta^{(m)}$, we encode it as the external current $\boldsymbol{\mu}_s^{(m)} = W_{\text{fb}}\mathbf{z}$ using Eq. (8), and set

$$\bar{\boldsymbol{\mu}} = gJ\boldsymbol{\mu} + W_{\text{fb}}\mathbf{z} + \boldsymbol{\mu}_s^{(m)} \tag{11}$$
$$\bar{C} = g^2 JCJ^\top + \boldsymbol{\sigma}_s^2 I \tag{12}$$

in Eq. (2). We then simulate the network until its state convergence to a fixed point, denoted as $(\boldsymbol{\mu}^{(m)}, C^{(m)})$. After collecting the fixed points for every training feature $\boldsymbol{\mu}^{(m)}$, we train $W_{\text{out}}$ by minimizing the loss function

$$\sum_{i=m}^{M} \left\| W_{\text{out}}^\top \boldsymbol{\mu}^{(m)} - \mathbf{z}^{(m)} \right\|^2 + \alpha \left\| W_{\text{out}} \right\|_2^2, \tag{13}$$

where $\alpha$ is the factor of regularization. Note that the covariance is not contained in the loss function. We set $M = 20$, $\sigma_s^2 = 1$, and $\alpha = 5 \times 10^{-3}$ throughout the work, except in Fig. 6 and Fig. 7, where we vary these parameters to investigate their respective effects on the performance of the MNN.

**Network inference.** During the cue period, to input a variable value $\theta^*$ for the model to hold as a memory item, we calculate the corresponding external input $\mathbf{z}^*$ and add external noise to the network by setting $\boldsymbol{\mu}_s$ as $W_{\text{fb}}\mathbf{z}^* + \boldsymbol{\xi}_t$, where $\boldsymbol{\xi}_t \sim \mathcal{N}(\mathbf{0}, \sigma_\xi^2 I)$ is a Gaussian noise term, with $\sigma_\xi \in \text{Unif}[0, 1]$ independently sampled for each instance. The mean $\bar{\boldsymbol{\mu}}$ and covariance $\bar{C}$ of the synaptic inputs in Eq. (2) as

$$\bar{\boldsymbol{\mu}} = (gJ + W_{\text{fb}}W_{\text{out}}^\top)\boldsymbol{\mu} + \boldsymbol{\mu}_s \tag{14}$$
$$\bar{C} = (gJ + W_{\text{fb}}W_{\text{out}}^\top)C(gJ + W_{\text{fb}}W_{\text{out}}^\top)^\top + \sigma_{s,2}^2 I, \tag{15}$$

where $\sigma_{s,2}$ is the noise level during the inference phase. The cue period lasts for 500 time steps, after that we remove the external inputs by setting $\boldsymbol{\mu}_s = 0$. Due the delay period, we simulate the network for 7500 steps. After the delay period, we decode the remembered feature from the network. The remembered variable is decoded from the mean as $\boldsymbol{\mu}_z = W_{\text{out}}^\top \boldsymbol{\mu}$, and the corresponding covariance is calculated as $\hat{C}_z = W_{\text{out}}^\top C W_{\text{out}}$. We calculate the error of the network as $e = |\hat{\theta} - \theta^*|$, where $\hat{\theta}$ and $\theta^*$ are the angles corresponds to $\hat{\boldsymbol{\mu}}_z$ and $\mathbf{z}^*$ respectively.

**Uncertainty metrics.** Based on the entropy of a Gaussian distribution, we design the following two uncertainty indicators

$$\text{(I): } \kappa = \det|\hat{C}_z|, \quad \text{(II): } \kappa = \log \det|\hat{C}_z|. \tag{16}$$

Denote the normal direction of the tangent space of $\mathcal{M}$ at $\hat{\boldsymbol{\mu}}_z$ as $n_{\hat{\boldsymbol{\mu}}_z} \in \mathbb{R}^2$, then we design another two uncertainty indicators

$$\text{(III): } \kappa = \det|n_{\hat{\boldsymbol{\mu}}_z}^\top \hat{C}_z n_{\hat{\boldsymbol{\mu}}_z}|, \quad \text{(IV): } \kappa = \log \det|n_{\hat{\boldsymbol{\mu}}_z}^\top \hat{C}_z n_{\hat{\boldsymbol{\mu}}_z}|. \tag{17}$$

These two uncertainty metrics measure the uncertainty along the tangent space of $\hat{\boldsymbol{\mu}}_z$ the task-related direction.

**Correlation between the uncertainty metrics and the error.** After trained a network, we select $L$ different angle samples $\theta_1^*, \ldots, \theta_L^*$ and generate the corresponding external stimuli, which are exerted to the network separately. For each $\theta_i^*$, we calculate the error $e_i$ and the uncertainty $u_i$ ($u_i$ is calculated through one of I-IV indicators). We then estimate the correlation

$$\rho = \frac{1}{L} \sum_i \kappa_i e_i - \frac{1}{L^2} \sum_i \kappa_i \sum_i e_i. \tag{18}$$

A large (positive) $\rho$ means that the uncertainty quantification is good. Throughout this work, we set $L = 500$ and uniformly select $\theta_1^*, \ldots, \theta_L^*$ from $[0, 2\pi)$.

**Differential covariance ratio (DCR).** We define the differential covariance ratio (DCR) as follows

$$\text{Diff}(C) := \frac{\sum_{ij}|C_{ij}|^2}{\sum_{ij}|C_{ij}|^2 + \sum_{ij}|C_{ij} + C_{i,2N-j} + C_{i,2N-j} + C_{2N-i,2N-j}|^2}. \tag{19}$$

The DCR quantifies the differential covariance within the input covariance matrix $C$, as the differential covariance induced by the drift of the bump exhibits a distinct pattern of sign relations (see Fig. 4)

$$\text{sgn}(C_{ij}) = -\text{sgn}(C_{2N-i,j}) = -\text{sgn}(C_{i,2N-j}) = \text{sgn}(C_{2N-i,2N-j}) = 1. \tag{20}$$

The higher the DCR, the larger the component of differential covariance in $C$. The maximum value of the DCR is 1.

**Simulation and decoding of the spiking neural network.** For Fig. 5, we simulated an SNN using the LIF neural model (Eq. (3)). The durations of the cue and delay periods were set to match those used in the MNN simulation. To decode the feature and its uncertainty from the spike trains, we select a time interval $[T_{\text{start}}, T_{\text{end}}]$, and divided it into $K = \frac{T_{\text{end}} - T_{\text{start}}}{\Delta T}$ time windows, where the window length $\Delta T$ was chosen such that it evenly divided $T_{\text{end}} - T_{\text{start}}$. For each time window indexed by $k$ ($k = 1, \ldots, K$), we calculated the mean firing rate of each neuron within the window, denoted as $\bar{\boldsymbol{r}}^{(k)} \in \mathbb{R}^N$. In this study, we set $T_{\text{start}} = 2000$, $T_{\text{end}} = 7500$ and $\Delta T = 500$. The mean firing rate and firing covariance of the SNN were then estimated as:

$$\boldsymbol{\mu}_{\text{snn}} = \frac{1}{K} \sum_{k=1}^K \bar{\boldsymbol{r}}^{(k)}, \quad C_{\text{snn}} = \frac{1}{K} \sum_{k=1}^K \left(\bar{\boldsymbol{r}}^{(k)} - \boldsymbol{\mu}_{\text{snn}}\right)\left(\bar{\boldsymbol{r}}^{(k)} - \boldsymbol{\mu}_{\text{snn}}\right)^\top.$$

The subsequent decoding procedure was identical to that used for the MNN. For details, refer to the **Encoding and decoding in the MNN** section.

# References

[1] Shaiyan Keshvari, Ronald Van den Berg, and Wei Ji Ma. Probabilistic computation in human perception under variability in encoding precision. *PLoS One*, 7(6):e40216, 2012.

[2] Deepna Devkar, Anthony A Wright, and Wei Ji Ma. Monkeys and humans take local uncertainty into account when localizing a change. *Journal of Vision*, 17(11):4–4, 2017.

[3] Maija Honig, Wei Ji Ma, and Daryl Fougnie. Humans incorporate trial-to-trial working memory uncertainty into rewarded decisions. *Proceedings of the National Academy of Sciences*, 117(15):8391–8397, 2020.

[4] Syaheed B Jabar, Kartik K Sreenivasan, Stergiani Lentzou, Anish Kanabar, Timothy F Brady, and Daryl Fougnie. Using a betting game to reveal the rich nature of visual working memories. *BioRxiv*, pages 2020–10, 2020.

[5] Aspen H Yoo, Luigi Acerbi, and Wei Ji Ma. Uncertainty is maintained and used in working memory. *Journal of vision*, 21(8):13–13, 2021.

[6] Jakob Gawlikowski, Cedrique Rovile Njieutcheu Tassi, Mohsin Ali, Jongseok Lee, Matthias Humt, Jianxiang Feng, Anna Kruspe, Rudolph Triebel, Peter Jung, Ribana Roscher, et al. A survey of uncertainty in deep neural networks. *Artificial Intelligence Review*, 56(Suppl 1):1513–1589, 2023.

[7] Edgar Y Walker, Stephan Pohl, Rachel N Denison, David L Barack, Jennifer Lee, Ned Block, Wei Ji Ma, and Florent Meyniel. Studying the neural representations of uncertainty. *Nature neuroscience*, 26(11):1857–1867, 2023.

[8] Dominik Seuß. Bridging the gap between explainable ai and uncertainty quantification to enhance trustability. *arXiv preprint*, 2021.

[9] Jochen Gast and Stefan Roth. Lightweight probabilistic deep networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3369–3378, 2018.

[10] Janis Postels, Francesco Ferroni, Huseyin Coskun, Nassir Navab, and Federico Tombari. Sampling-free epistemic uncertainty estimation using approximated variance propagation. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 2931–2940, 2019.

[11] Wei Ji Ma and Mehrdad Jazayeri. Neural coding of uncertainty and probability. *Annual review of neuroscience*, 37(1):205–220, 2014.

[12] Wei Ji Ma, Jeffrey M Beck, Peter E Latham, and Alexandre Pouget. Bayesian inference with probabilistic population codes. *Nature neuroscience*, 9(11):1432–1438, 2006.

[13] Gergő Orbán, Pietro Berkes, József Fiser, and Máté Lengyel. Neural variability and sampling-based probabilistic representations in the visual cortex. *Neuron*, 92(2):530–543, 2016.

[14] Hsin-Hung Li, Thomas C Sprague, Aspen H Yoo, Wei Ji Ma, and Clayton E Curtis. Joint representation of working memory and uncertainty in human cortex. *Neuron*, 109(22):3699–3712, 2021.

[15] Paul M Bays, Sebastian Schneegans, Wei Ji Ma, and Timothy F Brady. Representation and computation in visual working memory. *Nature Human Behaviour*, pages 1–19, 2024.

[16] Klaus Oberauer, Stephan Lewandowsky, Edward Awh, Gordon DA Brown, Andrew Conway, Nelson Cowan, Christopher Donkin, Simon Farrell, Graham J Hitch, Mark J Hurlstone, et al. Benchmarks for models of short-term and working memory. *Psychological bulletin*, 144(9):885, 2018.

[17] Ronald Van den Berg, Hongsup Shin, Wen-Chuang Chou, Ryan George, and Wei Ji Ma. Variability in encoding precision accounts for visual short-term memory limitations. *Proceedings of the National Academy of Sciences*, 109(22):8780–8785, 2012.

[18] Rosanne L Rademaker, Caroline H Tredway, and Frank Tong. Introspective judgments predict the precision and likelihood of successful maintenance of visual working memory. *Journal of vision*, 12(13):21–21, 2012.

[19] Olivier J Hénaff, Zoe M Boundy-Singer, Kristof Meding, Corey M Ziemba, and Robbe LT Goris. Representation of visual uncertainty through neural gain variability. *Nature communications*, 11(1):2513, 2020.

[20] Klaus Wimmer, Duane Q Nykamp, Christos Constantinidis, and Albert Compte. Bump attractor dynamics in prefrontal cortex explains behavioral precision in spatial working memory. *Nature neuroscience*, 17(3):431–439, 2014.

[21] Sung Soo Kim, Hervé Rouault, Shaul Druckmann, and Vivek Jayaraman. Ring attractor dynamics in the drosophila central brain. *Science*, 356(6340):849–853, 2017.

[22] Rishidev Chaudhuri, Berk Gerçek, Biraj Pandey, Adrien Peyrache, and Ila Fiete. The intrinsic attractor manifold and population dynamics of a canonical cognitive circuit across waking and sleep. *Nature neuroscience*, 22(9):1512–1520, 2019.

[23] Mikail Khona and Ila R Fiete. Attractor and integrator networks in the brain. *Nature Reviews Neuroscience*, 23(12):744–766, 2022.

[24] Anna Kutschireiter, Melanie A Basnak, Rachel I Wilson, and Jan Drugowitsch. Bayesian inference in ring attractor networks. *Proceedings of the National Academy of Sciences*, 120(9):e2210622120, 2023.

[25] Yoram Burak and Ila R Fiete. Accurate path integration in continuous attractor network models of grid cells. *PLoS computational biology*, 5(2):e1000291, 2009.

[26] Si Wu, KY Michael Wong, CC Alan Fung, Yuanyuan Mi, and Wenhao Zhang. Continuous attractor neural networks: candidate of a canonical model for neural information representation. *F1000Research*, 5, 2016.

[27] Dimitry Fisher, Itsaso Olasagasti, David W Tank, Emre RF Aksay, and Mark S Goldman. A modeling framework for deriving the structural and functional architecture of a short-term memory microcircuit. *Neuron*, 79(5):987–1000, 2013.

[28] Valerio Mante, David Sussillo, Krishna V Shenoy, and William T Newsome. Context-dependent computation by recurrent dynamics in prefrontal cortex. *nature*, 503(7474):78–84, 2013.

[29] H Francis Song, Guangyu R Yang, and Xiao-Jing Wang. Training excitatory-inhibitory recurrent neural networks for cognitive tasks: a simple and flexible framework. *PLoS computational biology*, 12(2):e1004792, 2016.

[30] A Emin Orhan and Wei Ji Ma. A diverse range of factors affect the nature of neural representations underlying short-term memory. *Nature neuroscience*, 22(2):275–283, 2019.

[31] Ran Darshan and Alexander Rivkind. Learning to represent continuous variables in heterogeneous neural networks. *Cell Reports*, 39(1), 2022.

[32] Rodrigo Echeveste, Laurence Aitchison, Guillaume Hennequin, and Máté Lengyel. Cortical-like dynamics in recurrent circuits optimized for sampling-based probabilistic inference. *Nature neuroscience*, 23(9):1138–1149, 2020.

[33] Wenlian Lu, Enrico Rossoni, and Jianfeng Feng. On a gaussian neuronal field model. *NeuroImage*, 52(3):913–933, 2010.

[34] Hengyuan Ma, Yang Qi, Pulin Gong, Jie Zhang, Wen-lian Lu, and Jianfeng Feng. Self-organization of nonlinearly coupled neural fluctuations into synergistic population codes. *Neural Computation*, 35(11):1820–1849, 2023.

[35] Zhichao Zhu, Yang Qi, Wenlian Lu, and Jianfeng Feng. Learning to integrate parts for whole through correlated neural variability. *PLOS Computational Biology*, 20(9):e1012401, 2024.

[36] Zhichao Zhu, Yang Qi, Wenlian Lu, Zhigang Wang, Lu Cao, and Jianfeng Feng. Towards free-response paradigm: a theory on decision-making in spiking neural networks. *arXiv preprint arXiv:2404.10599*, 2024.

[37] Klaus Oberauer. Access to information in working memory: exploring the focus of attention. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, 28(3):411, 2002.

[38] Yang Qi, Michael Breakspear, and Pulin Gong. Subdiffusive dynamics of bump attractors: mechanisms and functional roles. *Neural computation*, 27(2):255–280, 2015.

[39] Yoram Burak and Ila R Fiete. Fundamental limits on persistent activity in networks of noisy neurons. *Proceedings of the National Academy of Sciences*, 109(43):17645–17650, 2012.

[40] Daryl Fougnie, Jordan W Suchow, and George A Alvarez. Variability in the quality of visual working memory. *Nature communications*, 3(1):1229, 2012.

[41] Anne K Churchland, Roozbeh Kiani, Rishidev Chaudhuri, Xiao-Jing Wang, Alexandre Pouget, and Michael N Shadlen. Variance as a signature of neural computations during decision making. *Neuron*, 69(4):818–831, 2011.

[42] Amber Polk, Ashok Litwin-Kumar, and Brent Doiron. Correlated neural variability in persistent state networks. *Proceedings of the National Academy of Sciences*, 109(16):6295–6300, 2012.

[43] Moritz Helias, Tom Tetzlaff, and Markus Diesmann. The correlation structure of local neuronal networks intrinsically results from recurrent dynamics. *PLoS computational biology*, 10(1):e1003428, 2014.

[44] David Dahmen, Hannah Bos, and Moritz Helias. Correlated fluctuations in strongly coupled binary networks beyond equilibrium. *Physical Review X*, 6(3):031024, 2016.

[45] Shun-ichi Amari. Dynamics of pattern formation in lateral-inhibition type neural fields. *Biological cybernetics*, 27(2):77–87, 1977.

[46] Stephen Coombes. Waves, bumps, and patterns in neural field theories. *Biological cybernetics*, 93:91–108, 2005.

[47] Stephen Coombes and Helmut Schmidt. Neural fields with sigmoidal firing rates: approximate solutions. *Discret. Contin. Dyn. Syst. Ser. A*, 28:1369–1379, 2010.

[48] Daniel J Amit and MV Tsodyks. Quantitative study of attractor neural network retrieving at low spike rates. i. substrate-spikes, rates and neuronal gain. *Network: Computation in neural systems*, 2(3):259, 1991.

[49] Daniel J Amit and Nicolas Brunel. Dynamics of a recurrent network of spiking neurons before and following learning. *Network: Computation in Neural Systems*, 8(4):373–404, 1997.

[50] Alfonso Renart, Pengcheng Song, and Xiao-Jing Wang. Robust spatial working memory through homeostatic synaptic scaling in heterogeneous cortical networks. *Neuron*, 38(3):473–485, 2003.

[51] John D Murray, Alberto Bernacchia, Nicholas A Roy, Christos Constantinidis, Ranulfo Romo, and Xiao-Jing Wang. Stable population coding for working memory coexists with heterogeneous neural dynamics in prefrontal cortex. *Proceedings of the National Academy of Sciences*, 114(2):394–399, 2017.

[52] Christopher Langdon, Mikhail Genkin, and Tatiana A Engel. A unifying perspective on neural manifolds and circuits for cognition. *Nature Reviews Neuroscience*, 24(6):363–377, 2023.

[53] Zachary P Kilpatrick, Bard Ermentrout, and Brent Doiron. Optimizing working memory with heterogeneity of recurrent cortical excitation. *Journal of neuroscience*, 33(48):18999–19011, 2013.

[54] Charles Blundell, Julien Cornebise, Koray Kavukcuoglu, and Daan Wierstra. Weight uncertainty in neural network. In *International conference on machine learning*, pages 1613–1622. PMLR, 2015.

[55] Anqi Wu, Sebastian Nowozin, Edward Meeds, Richard E Turner, Jose Miguel Hernandez-Lobato, and Alexander L Gaunt. Deterministic variational inference for robust bayesian neural networks. In *International Conference on Learning Representations*, 2019.

[56] Wolfgang Maass, Thomas Natschläger, and Henry Markram. Real-time computing without stable states: A new framework for neural computation based on perturbations. *Neural computation*, 14(11):2531–2560, 2002.

[57] Michael Beyeler, Emily L Rounds, Kristofor D Carlson, Nikil Dutt, and Jeffrey L Krichmar. Neural correlates of sparse coding and dimensionality reduction. *PLoS computational biology*, 15(6):e1006908, 2019.

[58] Pietro Berkes, Gergő Orbán, Máté Lengyel, and József Fiser. Spontaneous cortical activity reveals hallmarks of an optimal internal model of the environment. *Science*, 331(6013):83–87, 2011.

[59] Jeffrey M Beck, Peter E Latham, and Alexandre Pouget. Marginalization in neural circuits with divisive normalization. *Journal of Neuroscience*, 31(43):15310–15319, 2011.

[60] Adam Kohn, Ruben Coen-Cagli, Ingmar Kanitscheider, and Alexandre Pouget. Correlations and neuronal population information. *Annual review of neuroscience*, 39(1):237–256, 2016.

[61] Stefano Panzeri, Monica Moroni, Houman Safaai, and Christopher D Harvey. The structures and functions of correlations in neural population codes. *Nature Reviews Neuroscience*, 23(9):551–567, 2022.

[62] Bruno B Averbeck, Peter E Latham, and Alexandre Pouget. Neural correlations, population coding and computation. *Nature reviews neuroscience*, 7(5):358–366, 2006.

[63] Laurence Aitchison, Jannes Jegminat, Jorge Aurelio Menendez, Jean-Pascal Pfister, Alexandre Pouget, and Peter E Latham. Synaptic plasticity as bayesian inference. *Nature neuroscience*, 24(4):565–571, 2021.

[64] Jianfeng Feng, Yingchun Deng, and Enrico Rossoni. Dynamics of moment neuronal networks. *Physical Review E—Statistical, Nonlinear, and Soft Matter Physics*, 73(4):041906, 2006.

[65] Yang Qi. Moment neural network and an efficient numerical method for modeling irregular spiking activity. *Phys. Rev. E*, 110:024310, Aug 2024.

[66] Jie Wu and Liqun Zhang. Backward uniqueness for general parabolic operators in the whole space. *Calculus of Variations and Partial Differential Equations*, 58:1–19, 2019.

[67] Lawrence C Evans. *Partial differential equations*, volume 19. American Mathematical Society, 2022.

# Supplementary Information: Uncertainty Quantification in Working Memory via Moment Neural Networks

Hengyuan Ma[1], Wenlian Lu[1,2,3,4,5,6], Jianfeng Feng[1,2,3,4,7*]

1 Institute of Science and Technology for Brain-inspired Intelligence, Fudan University, Shanghai 200433, China
2 Key Laboratory of Computational Neuroscience and Brain-Inspired Intelligence (Fudan University), Ministry of Education, China
3 School of Mathematical Sciences, Fudan University, No. 220 Handan Road, Shanghai, 200433, Shanghai, China
4 Shanghai Center for Mathematical Sciences, No. 220 Handan Road, Shanghai, 200433, Shanghai, China
5 Shanghai Key Laboratory for Contemporary Applied Mathematics, No. 220 Handan Road, Shanghai, 200433, Shanghai, China
6 Key Laboratory of Mathematics for Nonlinear Science, No. 220 Handan Road, Shanghai, 200433, Shanghai, China
7 Department of Computer Science, University of Warwick, Coventry, CV4 7AL, UK
∗ jffeng@fudan.edu.cn

## S1 Theorems and their proofs

Here, a theoretical analysis of an abstract neural system is presented, demonstrating how the system can learn to output covariance that effectively quantifies uncertainty by training under a loss function that explicitly supervises only the mean of the network output.

We represent the neural system as a simplified parameterized function $f(\mathbf{x}, \boldsymbol{\vartheta}) : \mathbb{R}^m \to \mathbb{R}$, where $\boldsymbol{\vartheta}$ denotes the learnable parameters, and $\mathbf{x}$ is the input. Let $f^*(\mathbf{x}) : \mathbb{R}^m \to \mathbb{R}$ represent the ground-truth neural system that the model aims to approximate. Internal and external noise are introduced to both systems: under specified noise conditions, the neuron's inference result is computed as

$$y(\mathbf{x}, C_\xi, \sigma_\eta^2; \boldsymbol{\vartheta}) = f(\mathbf{x} + \boldsymbol{\xi}_1; \boldsymbol{\vartheta}) + \eta_1, \tag{S1}$$

and the ground-truth result is

$$y^*(\mathbf{x}, C_\xi, \sigma_\eta^2) = f^*(\mathbf{x} + \boldsymbol{\xi}_2) + \eta_2, \tag{S2}$$

where $\boldsymbol{\xi}_1, \boldsymbol{\xi}_2 \sim \mathcal{N}(\mathbf{0}, C_\xi)$ is the internal noise with covariance matrix $C_\xi \in \mathbb{R}^{m \times m}$, and $\eta_1, \eta_2 \sim \mathcal{N}(0, \sigma_\eta^2)$ is the external noise with the variance $\sigma_\eta^2$. We suppose that $C_\xi$ is strictly positive-definite, and $\sigma_\eta^2 > 0$. $\boldsymbol{\xi}_1, \boldsymbol{\xi}_2, \eta_1, \eta_2$ are independent from each other. The model output mean and output variance are calculated as

$$m(\mathbf{x}, C_\xi, ; \boldsymbol{\vartheta}) = \mathbb{E}[y(\mathbf{x}, C_\xi, \sigma_\eta^2; \boldsymbol{\vartheta})], \tag{S3}$$

$$v(\mathbf{x}, C_\xi, \sigma_\eta^2; \boldsymbol{\vartheta}) = \mathbb{E}[\left(y(\mathbf{x}, C_\xi, \sigma_\eta^2; \boldsymbol{\vartheta}) - m(\mathbf{x}, C_\xi; \boldsymbol{\vartheta})\right)^2], \tag{S4}$$

and the ground-truth mean and variance is

$$m^*(\mathbf{x}, C_\xi) = \mathbb{E}[y^*(\mathbf{x}, C_\xi, \sigma_\eta^2)], \tag{S5}$$

$$v^*(\mathbf{x}, C_\xi, \sigma_\eta^2) = \mathbb{E}[\left(y^*(\mathbf{x}, C_\xi, \sigma_\eta^2) - m(\mathbf{x}, C_\xi)\right)^2]. \tag{S6}$$

Note that both $m^*$ and $m$ are not influenced by $\sigma_\eta^2$ since $\eta$ vanishes after taking the expectation, while both $m^*$ and $m$ are influenced by $C_\xi$ due to the nonlinearity of $f^*$ and $f$, respectively. The neural system is trained to minimize the following loss, which only supervises the mean under the noise conditions $C_\xi$ and $\sigma_\eta^2$ during training

$$\mathcal{L}(\boldsymbol{\vartheta}, C_\xi) = \int_{\mathbf{x} \in \mathbb{R}^m} \left(m^*(\mathbf{x}, C_\xi; \boldsymbol{\vartheta}) - m(\mathbf{x}, C_\xi; \boldsymbol{\vartheta})\right)^2 p(\mathbf{x}) d\mathbf{x}, \tag{S7}$$

where $p(\mathbf{x})$ is the distribution of the input, which is assumed to be supported on the whole $\mathbb{R}^m$.

We first prove a theorem demonstrating that the neural system can learn the ground-truth output variance by minimizing the loss $\mathcal{L}(\boldsymbol{\vartheta}, C_\xi)$, which supervises only the mean output.

**Theorem S1.** *Supposed that the model parameter $\boldsymbol{\vartheta}$ diminishes the loss $\mathcal{L}(\boldsymbol{\vartheta}, C_\xi)$, then the model learns the ground-truth output mean and variance at the same time for any noise conditions with positive-definite $C_\xi'$, and $\sigma_\eta'^2$*

$$m(\mathbf{x}, C_\xi'; \boldsymbol{\vartheta}) = m^*(\mathbf{x}, C_\xi'), \tag{S8}$$

$$v(\mathbf{x}, C_\xi', \sigma_\eta'^2; \boldsymbol{\vartheta}) = v^*(\mathbf{x}, C_\xi', \sigma_\eta'^2). \tag{S9}$$

*Proof.* We note that the function $m(\mathbf{x}, tC_\xi; \boldsymbol{\vartheta})$ is the solution to the following parabolic differential equation

$$\begin{cases} \partial_t w(t, \mathbf{x}) = \nabla_{\mathbf{x}}[C_{\boldsymbol{\xi}} \nabla_{\mathbf{x}} w(t, \mathbf{x})], & t > 0, \quad \mathbf{x} \in \mathbb{R}^n \\ w(t, \mathbf{x}) = m(\mathbf{x}, O; \boldsymbol{\vartheta}), & t = 0, \quad \mathbf{x} \in \mathbb{R}^n \end{cases}, \tag{S10}$$

where $O \in \mathbb{R}^{m \times m}$ is the all-zero matrix. Additionally, the function $m^*(\mathbf{x}, tC_\xi)$ is the solution to the following the same parabolic differential equation with a different initial condition

$$\begin{cases} \partial_t w(t, \mathbf{x}) = \nabla_{\mathbf{x}}[C_{\boldsymbol{\xi}} \nabla_{\mathbf{x}} w(t, \mathbf{x})], & t > 0, \quad \mathbf{x} \in \mathbb{R}^n \\ w(t, \mathbf{x}) = m^*(\mathbf{x}, O), & t = 0, \quad \mathbf{x} \in \mathbb{R}^n \end{cases}, \tag{S11}$$

Since $\mathcal{L}(\boldsymbol{\vartheta}, C_\xi, \sigma_\eta^2) = 0$, and $p(\mathbf{x})$ is support on the $\mathbb{R}^m$, we have

$$m(\mathbf{x}, C_\xi; \boldsymbol{\vartheta}) = m^*(\mathbf{x}, C_\xi), \quad \forall \mathbf{x} \in \mathbb{R}^m. \tag{S12}$$

This means that the solutions to the two parabolic differential equations, Eq. (S10) and Eq. (S11), are identical at $t = 1$. By applying the backward uniqueness property of parabolic differential equations [66], which states that if the solutions at a given time point $t > 0$ are identical, then their initial states must also be identical, we have

$$m(\mathbf{x}, O; \boldsymbol{\vartheta}) = m^*(\mathbf{x}, O), \quad \forall \mathbf{x} \in \mathbb{R}^m, \tag{S13}$$

which is equivalent to

$$f(\mathbf{x}; \boldsymbol{\vartheta}) = f^*(\mathbf{x}), \quad \forall \mathbf{x} \in \mathbb{R}^m. \tag{S14}$$

Then both Eq. (S8) and Eq. (S9) hold for any $C_\xi', \sigma_\eta'^2$. $\square$

Importantly, this theorem also suggests that the neural system inherently generalizes its ability to quantify uncertainty across different noise conditions, characterized by variations in $C_\xi$ and $\sigma_\eta^2$. This generalization ability has been observed in Figs. 7d-f in the main paper

In general, minimizing the training loss to zero is impossible. However, we present a stronger theoretical result showing that the error in covariance can be effectively bounded by the error in the mean.

**Theorem S2.** *Supposed that both $f(\mathbf{x}; \boldsymbol{\vartheta})$ and $f(\mathbf{x})$ are bounded by $B > 0$. Denote the loss under noise level $tI, \sigma_\eta^2$ as*

$$\mathcal{L}(t; \boldsymbol{\vartheta}) = \mathcal{L}(\boldsymbol{\vartheta}, tI, \sigma_\eta^2). \tag{S15}$$

*Define the error on the variance as*

$$\mathcal{L}_v(t; \boldsymbol{\vartheta}) = \int_{\mathbf{x}} |v^*(\mathbf{x}, tI, \sigma_\eta^2) - v(\mathbf{x}, tI, \sigma_\eta^2; \boldsymbol{\vartheta})| p(\mathbf{x}) d\mathbf{x} \tag{S16}$$

*Then for any $\epsilon > 0$, there exists $t_2 > t_1 > 0$, such that*

$$\mathcal{L}_v(t_1; \boldsymbol{\vartheta}) \leq 4B\left(\epsilon + \sqrt{2}\sigma_\eta + 2\left(\mathcal{L}(t_2; \boldsymbol{\vartheta}) + \frac{m}{2} \int_{t_1}^{t_2} \frac{\mathcal{L}(t; \boldsymbol{\vartheta})}{t} dt\right)^{1/2}\right). \tag{S17}$$

*Proof.* First, we apply the parabolic differential equation (Eq. (S10) and Eq. (S11)) and the integration by parts, we have

$$\frac{d}{dt}\mathcal{L}(t; \boldsymbol{\vartheta}) = -2 \int_{\mathbb{R}^n} \left\| \nabla_{\mathbf{x}}\left(m^*(\mathbf{x}, C_\xi; \boldsymbol{\vartheta}) - y^*(\mathbf{x}, C_\xi; \boldsymbol{\vartheta})\right) \right\|^2 p(\mathbf{x}) d\mathbf{x}. \tag{S18}$$

Hence we have for $0 < t_1 < t_2$

$$\mathcal{L}(t_1; \boldsymbol{\vartheta}) = \mathcal{L}(t_2; \boldsymbol{\vartheta}) + 2 \int_{t_1}^{t_2} \int_{\mathbb{R}^n} \left\| \nabla_{\mathbf{x}}\left(m^*(\mathbf{x}, C_\xi; \boldsymbol{\vartheta}) - y^*(\mathbf{x}, C_\xi; \boldsymbol{\vartheta})\right) \right\|^2 p(\mathbf{x}) d\mathbf{x} dt. \tag{S19}$$

Use the Harnack's inequality [67], we have

$$\left\| \nabla_{\mathbf{x}}\left(m^*(\mathbf{x}, C_\xi; \boldsymbol{\vartheta}) - y^*(\mathbf{x}, C_\xi; \boldsymbol{\vartheta})\right) \right\|^2 \leq r(\tau, \mathbf{x}; \boldsymbol{\theta}) \partial_\tau \left(m^*(\mathbf{x}, C_\xi; \boldsymbol{\vartheta}) - y^*(\mathbf{x}, C_\xi; \boldsymbol{\vartheta})\right) \tag{S20}$$

$$+ \frac{m}{2\tau}\left(m^*(\mathbf{x}, C_\xi; \boldsymbol{\vartheta}) - y^*(\mathbf{x}, C_\xi; \boldsymbol{\vartheta})\right)^2, \tag{S21}$$

2

combine with Eq. (S19), we have

$$\mathcal{L}(t_1; \boldsymbol{\vartheta}) \leq \mathcal{L}(t_2; \boldsymbol{\vartheta}) + \frac{m}{2} \int_{t_1}^{t_2} \frac{\mathcal{L}(t; \boldsymbol{\vartheta})}{t} dt. \tag{S22}$$

Using the Hölder inequality and Minkowski inequality on the measure $p(\mathbf{x})$, we have

$$\int_{\mathbf{x}} |y^*(\mathbf{x}, \boldsymbol{\vartheta}, O, \sigma_\eta^2) - y(\mathbf{x}, \boldsymbol{\vartheta}, O, \sigma_\eta^2)| p(\mathbf{x}) d\mathbf{x} \leq \Big( \int_{\mathbf{x}} |y^*(\mathbf{x}, \boldsymbol{\vartheta}, O, \sigma_\eta^2) - y(\mathbf{x}, \boldsymbol{\vartheta}, O, \sigma_\eta^2)|^2 p(\mathbf{x}) d\mathbf{x} \Big)^{1/2}$$

$$\leq \Big( \int_{\mathbb{R}^n} (m(\mathbf{x}, \boldsymbol{\vartheta}, O) - m(\mathbf{x}, \boldsymbol{\vartheta}, t_1 I))^2 p(\mathbf{x}) d\mathbf{x} \Big)^{1/2}$$

$$+ \Big( \int_{\mathbb{R}^n} (m(\mathbf{x}, \boldsymbol{\vartheta}, t_1 I) - m^*(\mathbf{x}, t_1 I))^2 p(\mathbf{x}) d\mathbf{x} \Big)^{1/2}$$

$$+ \Big( \int_{\mathbb{R}^n} (m^*(\mathbf{x}, t_1 I) - m^*(\mathbf{x}, \mathbf{x}, O))^2 d\mathbf{x} \Big)^{1/2} + \sqrt{2}\sigma_\eta$$

$$= \Big( \int_{\mathbb{R}^n} (m(\mathbf{x}, \boldsymbol{\vartheta}, O) - m(\mathbf{x}, \boldsymbol{\vartheta}, t_1 I))^2 p(\mathbf{x}) d\mathbf{x} \Big)^{1/2} +$$

$$+ \Big( \int_{\mathbb{R}^n} (m^*(\mathbf{x}, t_1 I, \sigma_\eta^2) - m^*(\mathbf{x}, O))^2 d\mathbf{x} \Big)^{1/2} + \mathcal{L}(t_1; \boldsymbol{\vartheta})^{1/2} + \sqrt{2}\sigma_\eta$$

$$\tag{S23}$$

where $O \in \mathbb{R}^{m \times m}$ is the all-zero matrix. Noticed that for $a, b, a', b' \in \mathbb{R}$, we have

$$|(a - b)^2 - (a' - b')^2| \leq |a + a' - 2b||a - a'| + |2a' - b - b'||b - b'| \leq 4 \max\{|a|, |a'|, |b|, |b'|\}(|a - a'| + |b - b'|). \tag{S24}$$

we have

$$\int_{\mathbf{x}} |v^*(\mathbf{x}, t_1 I, \sigma_\eta^2) - v(\mathbf{x}, t_1 I, \sigma_\eta^2; \boldsymbol{\vartheta})| p(\mathbf{x}) d\mathbf{x} \tag{S25}$$

$$\leq 4B \left( \int_{\mathbf{x}} \mathbb{E}[|y^*(\mathbf{x}, t_1 I, \sigma_\eta^2) - y(\mathbf{x}, t_1 I, \sigma_\eta^2)|] p(\mathbf{x}) d\mathbf{x} + \int_{\mathbf{x}} |m^*(\mathbf{x}, t_1 I, \sigma_\eta^2) - m(\mathbf{x}, t_1 I, \sigma_\eta^2)| p(\mathbf{x}) d\mathbf{x} \right) \tag{S26}$$

$$\leq 4B \left( \left( \int_{\mathbf{x}} \mathbb{E}[(y^*(\mathbf{x}, t_1 I, \sigma_\eta^2) - y(\mathbf{x}, t_1 I, \sigma_\eta^2))^2] p(\mathbf{x}) d\mathbf{x} \right)^{1/2} + \mathcal{L}(t_1; \boldsymbol{\vartheta})^{1/2} \right) \tag{S27}$$

$$= 4B \left( \left( \int_{\mathbf{x}} (y^*(\mathbf{x}, O, \sigma_\eta^2) - y(\mathbf{x}, O, \sigma_\eta^2))^2 \tilde{p}(\mathbf{x}) d\mathbf{x} \right)^{1/2} + \mathcal{L}(t_1; \boldsymbol{\vartheta})^{1/2} \right), \tag{S28}$$

where $\tilde{p}(\mathbf{x})$ is the distribution of $\mathbf{x} + \boldsymbol{\xi}$ with $\mathbf{x} \sim p(\mathbf{x})$ and $\boldsymbol{\xi} \sim \mathcal{N}(\mathbf{0}, t_1 I)$. Due the continuity of the Gaussian convolution kernel and the assumption that $y$ and $y^*$ are bounded, for a given $\epsilon > 0$, there exists $t_1 > 0$ small enough, such that

$$\Big( \int_{\mathbb{R}^n} (m(\mathbf{x}, \boldsymbol{\vartheta}, O) - m(\mathbf{x}, \boldsymbol{\vartheta}, t_1 I))^2 p(\mathbf{x}) d\mathbf{x} \Big)^{1/2} + \Big( \int_{\mathbb{R}^n} (m^*(\mathbf{x}, t_1 I, \sigma_\eta^2) - m^*(\mathbf{x}, O))^2 d\mathbf{x} \Big)^{1/2} \tag{S29}$$

$$+ \left| \left( \int_{\mathbf{x}} (y^*(\mathbf{x}, O, \sigma_\eta^2) - y(\mathbf{x}, O, \sigma_\eta^2))^2 \tilde{p}(\mathbf{x}) d\mathbf{x} \right)^{1/2} - \left( \int_{\mathbf{x}} (y^*(\mathbf{x}, O, \sigma_\eta^2) - y(\mathbf{x}, O, \sigma_\eta^2))^2 p(\mathbf{x}) d\mathbf{x} \right)^{1/2} \right| \leq \epsilon. \tag{S30}$$

Then we have

$$\int_{\mathbf{x}} |v^*(\mathbf{x}, t_1 I, \sigma_\eta^2) - v(\mathbf{x}, t_1 I, \sigma_\eta^2; \boldsymbol{\vartheta})| p(\mathbf{x}) d\mathbf{x} \leq 4B(\epsilon + \sqrt{2}\sigma_\eta + 2\mathcal{L}(t_1; \boldsymbol{\vartheta})^{1/2}). \tag{S31}$$

Combine with Eq. (S22), we have

$$\int_{\mathbf{x}} |v^*(\mathbf{x}, t_1 I, \sigma_\eta^2) - v(\mathbf{x}, t_1 I, \sigma_\eta^2; \boldsymbol{\vartheta})| p(\mathbf{x}) d\mathbf{x} \leq 4B(\epsilon + \sqrt{2}\sigma_\eta + 2(\mathcal{L}(t_2; \boldsymbol{\vartheta}) + \frac{m}{2} \int_{t_1}^{t_2} \frac{\mathcal{L}(t; \boldsymbol{\vartheta})}{t} dt)^{1/2}), \tag{S32}$$

which proves the theorem. $\qquad \square$

This theorem states that the error in the output variance can be controlled by the mean loss under higher noise conditions ($t_2 > t_1$). This suggests that the neural system should be trained under a higher noise level to effectively quantify uncertainty within the noise range present during inference.

Both Thm. S1 and Thm. S2 can be extended to cases where the output $y$ is high-dimensional.
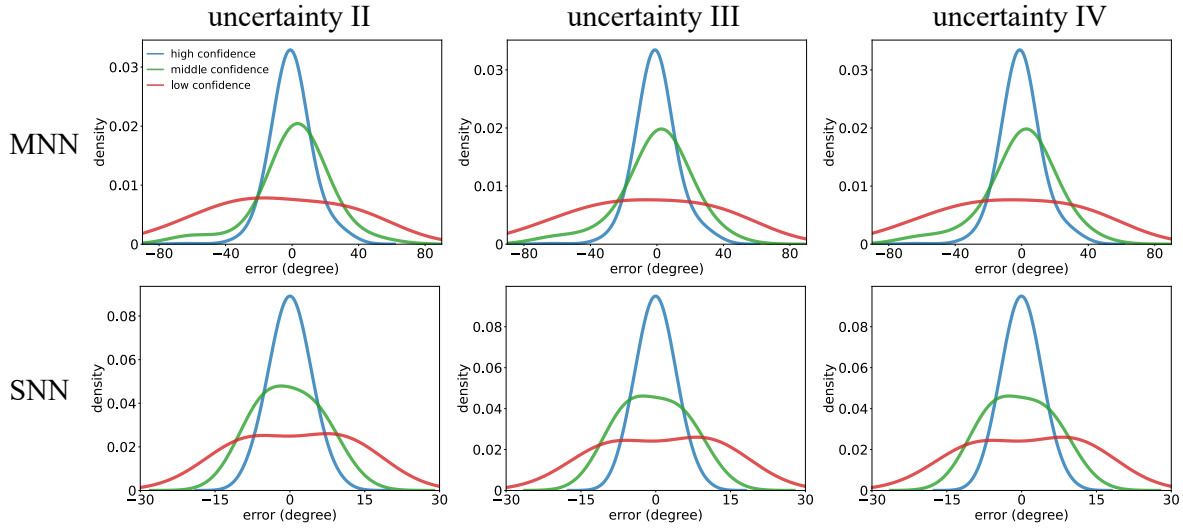
# S2   Supplementary figures



Figure S1: The error distribution across three groups of instances, divided based on the level of uncertainty (using the uncertainty metrics II, III, and IV (Eq. (16)-(17), *Methods*)): top 25% uncertainty (low confidence), top 25-50% uncertainty (middle confidence), and the remaining instances (high confidence). We plot the result for both moment neural network (MNN) and spiking neural network (SNN). The corresponding results calculated by the uncertainty metrics I are shown in Fig. 3d and Fig. 5c, respectively.

Figure S2: Effect of training samples $M$, population size $N$, and regularization $\alpha$ on the uncertainty quantification performance (evaluated by the correlation between the error and the uncertainty quantification calculated by the uncertainty metrics (Eq. (16)-(17), *Methods*). The corresponding results of uncertainty metrics I is shown in Fig. 6a,d,g.
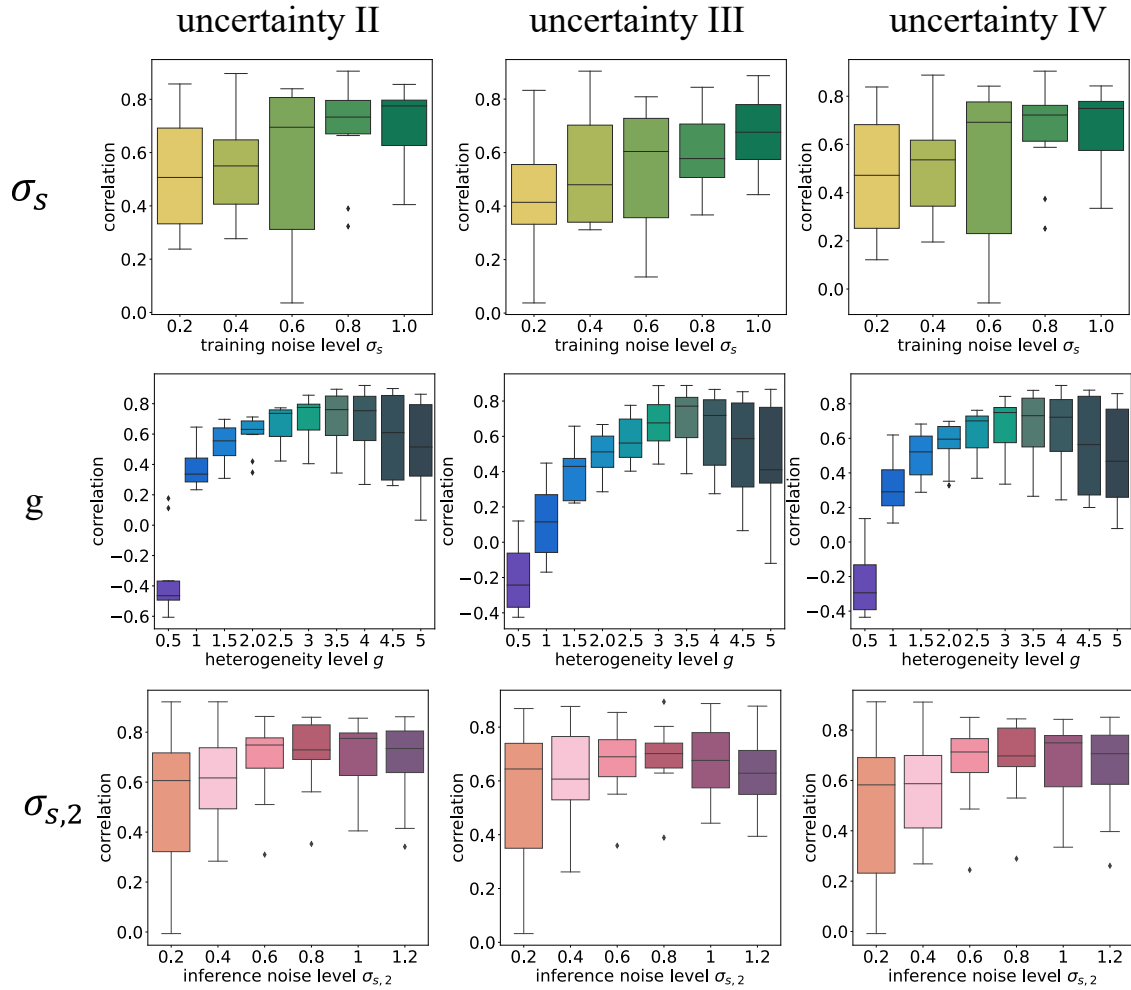
Figure S3: Effect of level of noise at training phase $\sigma_s$, heterogeneity $g$ and level of noise at inference phase $\sigma_{s,2}$ on the uncertainty quantification performance (evaluated by the correlation between the error and the uncertainty quantification calculated by the uncertainty metrics (Eq. (16)-(17), *Methods*)) The corresponding results of uncertainty metrics I is shown in Fig. 7a,d,g.
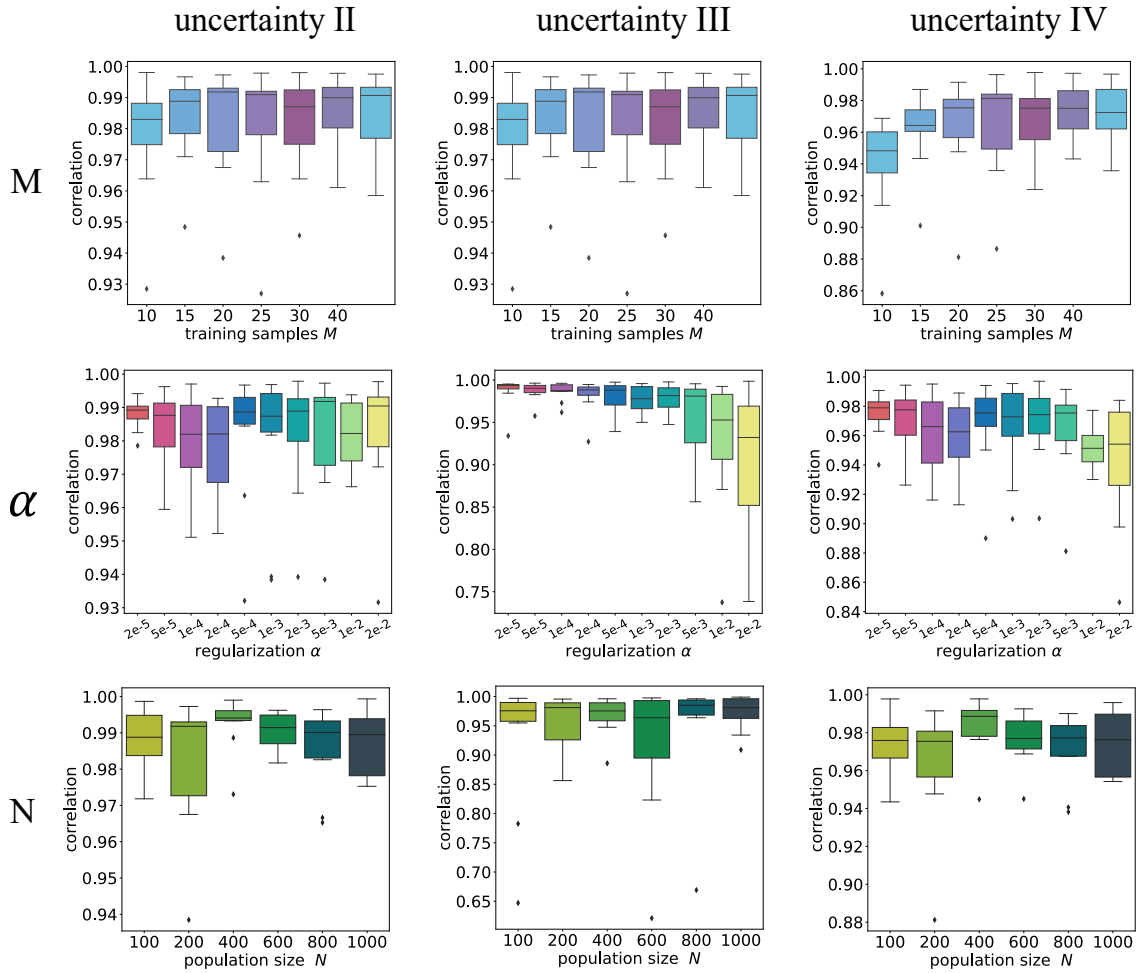
Figure S4: Effect of training samples $M$, population size $N$, and regularization $\alpha$ on the mean-covariance coupling strength (evaluated by the correlation between the bump width and the uncertainty quantification calculated by the uncertainty metrics (Eq. (16)-(17), *Methods*). The corresponding results of uncertainty metrics I is shown in Fig. 6b,e,h.
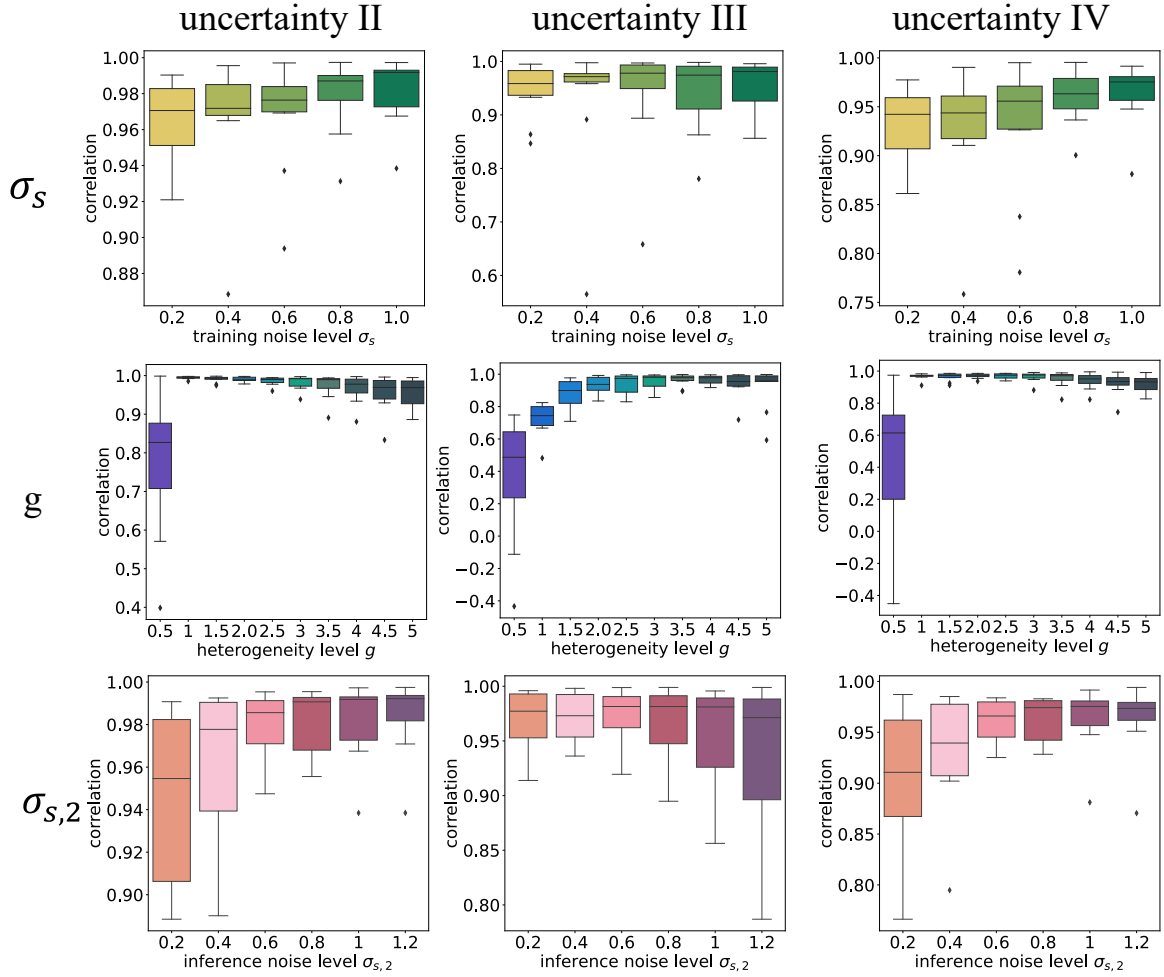
Figure S5: Effect of level of noise at training phase $\sigma_s$, heterogeneity $g$ and level of noise at inference phase $\sigma_{s,2}$ on the mean-covariance coupling strength (evaluated by the correlation between the bump width and the uncertainty quantification calculated by the uncertainty metrics (Eq. (16)-(17), *Methods*). The corresponding results of uncertainty metrics I is shown in Fig. 7b,e,h.