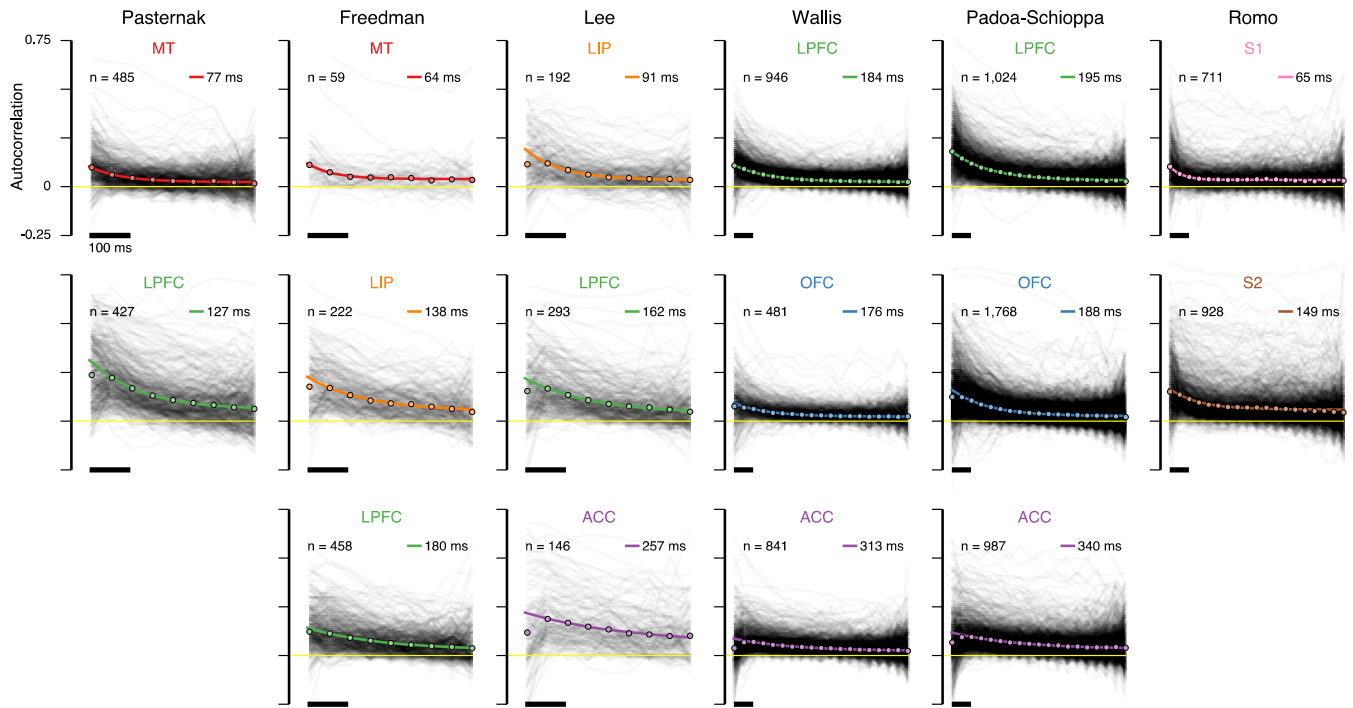**Supplementary Figure 1**

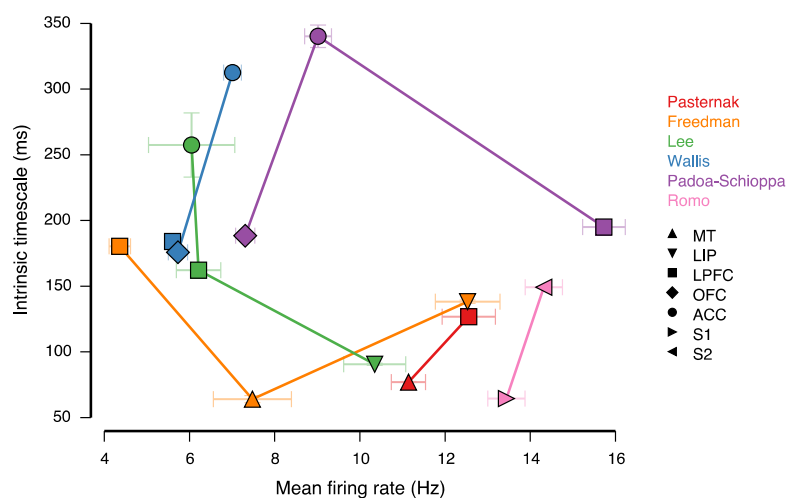Spike-count autocorrelations in time.

Normalized autocorrelation matrices are shown for each area in a dataset. The matrix shows the mean correlation of the spike count in each time bin with the spike count in every other time bin, averaged across neurons. These show that the autocorrelation is roughly stationary across time during the foreperiod.

**Supplementary Figure 2**

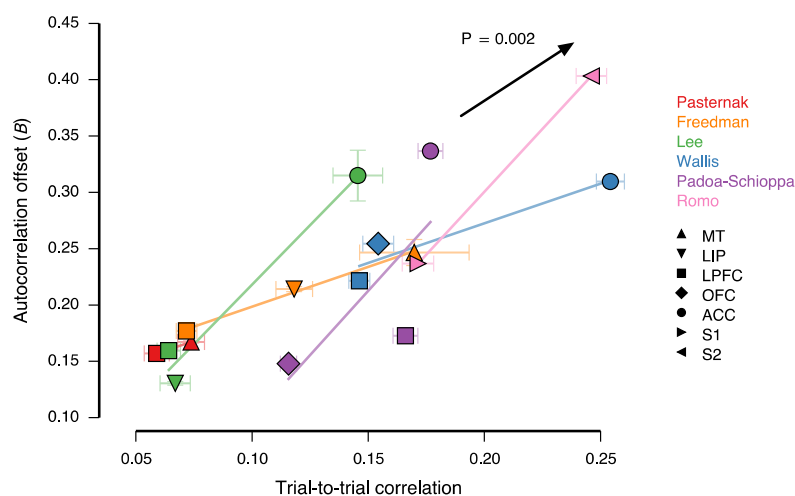Single neurons exhibit heterogeneous autocorrelations.

Light grey traces show the spike-count autocorrelation as function of time lag for single neurons, averaged across time points. Circles mark the population mean at each time lag, and the curve shows the exponential fit to the population data. The observation of single-neuron heterogeneity reinforces the interpretation of intrinsic timescale as a characteristic at the population level rather than at the single-neuron level.

**Supplementary Figure 3**

Differences in mean firing rates across areas do not account for hierarchy of intrinsic timescales.
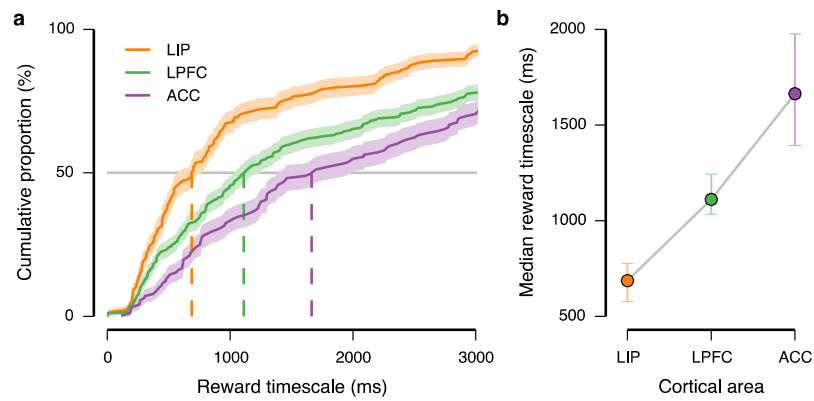
Mean firing rates varied substantially across datasets and across areas within datasets. There was no significant dependence of intrinsic timescale on mean firing rate ($P = 0.51$, $t(9) = -0.69$, two-tailed $t$-test, regression slope $m = -5.5 \pm 7.9$ ms/Hz; $P = 0.16$, $r_s = -0.34$, Spearman's rank correlation, two-tailed). Error bars mark s.e.

**Supplementary Figure 4**

Autocorrelation offset reflects trial-to-trial correlation.

Trial-to-trial correlation was calculated as the Pearson correlation coefficient between the foreperiod spike count in each trial and the spike count in the next trial. We hypothesized that autocorrelation offset would positively correlate with trial-to-trial correlation, and found a significant positive correlation between them. This indicates that the autocorrelation offset includes contributions from variability at timescales are comparable to or longer than the trial duration. Colored lines show trends for individual datasets. The arrow shows the slope of dependence from a regression analysis (slope $m = 1.3 \pm 0.3$). Error bars mark s.e.

**Supplementary Figure 5**

Hierarchical ordering of areas by timescale of reward memory.

In the Lee dataset, we previously measured timescales of the decay of memory traces for past rewards in single-neuron firing rates, while monkeys performed a competitive decision-making task. **(a)** The cumulative distribution of reward timescales in LIP (n = 160), LPFC (n = 243), and ACC (n = 134). For neurons fit with the sum of two reward timescales, we used the harmonic mean of the two timescales. **(b)** Median reward timescale for the three areas. Error bars mark s.e.

# Supplementary Mathematical Note

The mathematical framework of doubly stochastic point processes informs our interpretation of spike-count autocorrelation. Here we present calculations for the variability and autocorrelation for an inhomogeneous Poisson process with stochastic rate. We consider how rate fluctuations at different timescales affects the spike-count autocorrelation. Analysis of multiple timescales results in the functional form used to fit the autocorrelation data (Eq. S.20).

## Definitions

We consider a point process that generates a sequence of events at particular times. The count in a given time window $(t, t + \Delta)$ is denoted $N(t, \Delta)$. For a homogeneous Poisson process, the event rate $\lambda(t)$ is given by a constant $\mu$, and the mean and variance of the count are both equal to $\mu\Delta$. For an inhomogeneous Poisson process, the rate $\lambda(t)$ changes in time. In this case, the mean and variance are also equal, but calculated by the integral of the event rate:

$$\langle N(t, \Delta) \rangle = \int_t^{t+\Delta} ds \lambda(s), \text{ and} \tag{S.1}$$

$$\langle N(t, \Delta)^2 \rangle - \langle N(t, \Delta) \rangle^2 = \int_t^{t+\Delta} ds \lambda(s) \tag{S.2}$$

If $\lambda(t) = \mu$, this is the homogeneous Poisson process.

Here we are interested the case of a doubly stochastic process, in which the time-dependent rate $\lambda(t)$ is itself a stochastic process, characterized by a mean and variance. For simplicity, we assume the mean and variance of the event rate are constant in time, i.e. $\langle \lambda(t) \rangle = \mu$ and $\langle \lambda(t)^2 \rangle - \langle \lambda(t) \rangle^2 = \sigma^2$. Total variability is due both to variability in event timing and variability in event rate.

## Mean and variance

Averaging over both sources of stochasticity, the total mean of the count is equal to:

$$\langle \langle N(t, \Delta) \rangle \rangle = \int_t^{t+\Delta} ds \langle \lambda(s) \rangle = \mu\Delta \tag{S.3}$$

The mean count is equal to that of a Poisson process with constant rate $\mu$. The variance of counts will differ from that of a Poisson process, to reflect the contribution of rate variability. The second

1

moment of the count is given by:

$$\left\langle\left\langle N(t,\Delta)^2\right\rangle\right\rangle = \left\langle\int_t^{t+\Delta} ds\lambda(s) + \left[\int_t^{t+\Delta} ds\left\langle\lambda(s)\right\rangle\right]^2\right\rangle = \mu\Delta + \int_t^{t+\Delta} ds\int_t^{t+\Delta} ds'\left\langle\lambda(s)\lambda(s')\right\rangle \quad \text{(S.4)}$$

We denote $\delta\lambda(t) \equiv \lambda(t) - \mu$, the deviation of the rate from its mean. The total variance is then given by:

$$\left\langle\left\langle N(t,\Delta)^2\right\rangle\right\rangle - \left\langle\left\langle N(t,\Delta)\right\rangle\right\rangle^2 = \mu\Delta + \int_t^{t+\Delta} ds\int_t^{t+\Delta} ds'\left\langle\delta\lambda(s)\delta\lambda(s')\right\rangle \quad \text{(S.5)}$$

Here we consider that the event rate is characterized by a given time constant $\tau$. We define the autocovariance of the rate to be:

$$\left\langle\delta\lambda(s)\delta\lambda(s')\right\rangle = \sigma^2\exp\left(-\frac{|s-s'|}{\tau}\right) \quad \text{(S.6)}$$

For instance, an Orstein-Uhlenbeck process has this form of an autocovariance. Note that the linear rate model, shown in Fig. 4a and described in the Methods, behaves as an Ornstein-Uhlenbeck process in response to white-noise input and produces a rate autocovariance of the form in Eq. S.6. The variance of the count is then given by:

$$\left\langle\left\langle N(t,\Delta)^2\right\rangle\right\rangle - \left\langle\left\langle N(t,\Delta)\right\rangle\right\rangle^2 = \mu\Delta + 2\sigma^2\tau^2\left[\exp\left(-\frac{\Delta}{\tau}\right) - \left(1 - \frac{\Delta}{\tau}\right)\right] \quad \text{(S.7)}$$

The first term is the contribution of event timing variability and the second term is the contribution of event rate variability. If the time window $\Delta$ is small relative to $\tau$ ($\Delta \ll \tau$), the second term is approximately equal to $\sigma^2\Delta^2$ and independent of $\tau$. If the window size $\Delta$ is large relative to $\tau$ ($\Delta \gg \tau$), the second term is approximately equal to $2\sigma^2\tau\Delta$.

## Autocorrelation

We now examine how rate variability affects correlations in time. We consider covariance between counts in disjoint time bins, each of same window size $\Delta$, separated by an integer number of windows $k$ for a total time lag of $k\Delta$. The autocovariance $C$ at $k$-lag is defined as:

$$C(t, t + k\Delta) \equiv \left\langle\left\langle N(t,\Delta)N(t + k\Delta,\Delta)\right\rangle\right\rangle - \left\langle\left\langle N(t,\Delta)\right\rangle\right\rangle\left\langle\left\langle N(t + k\Delta,\Delta)\right\rangle\right\rangle \quad \text{(S.8)}$$

2

The event timings are independent because the intervals are disjoint ($k \geq 1$), so event timing variability does not contribute to autocovariance. The autocovariance is therefore:

$$C(k\Delta) = \int_t^{t+\Delta} ds \int_t^{t+\Delta} ds' \langle \delta\lambda(s)\delta\lambda(s') \rangle = \sigma^2 \left[ 2\tau \sinh\left(\frac{\Delta}{2\tau}\right) \exp\left(-\frac{k\Delta}{\tau}\right) \right]^2 \qquad \text{(S.9)}$$

When $k = 0$, this equation does not apply and autocovariance is simply equal to the variance. In the limit of slow rate fluctuations ($\Delta \ll \tau$):

$$C(k\Delta) \simeq \sigma^2 \Delta^2 \exp\left(-\frac{k\Delta}{\tau}\right) \qquad \text{(S.10)}$$

In the limit of fast rate fluctuations ($\Delta \gg \tau$):

$$C(k\Delta) \simeq \delta_{1k} \sigma^2 \tau^2 \qquad \text{(S.11)}$$

where $\delta_{ij}$ is the Kronecker delta function. For non-contiguous intervals ($k > 1$), $C \simeq 0$, and for contiguous intervals ($k = 1$), $C \simeq \sigma^2 \tau^2$.

The autocorrelation $R$ is given by the autocovariance divided by the variance:

$$R(k\Delta) = \left( \frac{\sigma^2 \left[ 2\tau \sinh\left(\frac{\Delta}{2\tau}\right) \right]^2}{\mu\Delta + 2\sigma^2 \tau^2 \left[ \exp\left(-\frac{\Delta}{\tau}\right) - \left(1 - \frac{\Delta}{\tau}\right) \right]} \right) \exp\left(-\frac{k\Delta}{\tau}\right) \qquad \text{(S.12)}$$

When $k = 0$, this equation does not apply and $R = 1$ by definition. In the limit of slow rate fluctuations ($\Delta \ll \tau$):

$$R(k\Delta) \simeq \left( 1 + \frac{\mu}{\sigma^2 \Delta} \right)^{-1} \exp\left(-\frac{k\Delta}{\tau}\right) \qquad \text{(S.13)}$$

If rate fluctuations are large ($\sigma^2 \Delta/\mu \gg 1$) the first factor is $\simeq 1$. If rate fluctuations are small ($\sigma^2 \Delta/\mu \ll 1$), this first factor is also small. In the limit of fast rate fluctuations ($\Delta \gg \tau$):

$$R(k\Delta) \simeq \delta_{1k} \left( \frac{\tau}{2\Delta} \right) \left( 1 + \frac{\mu}{2\sigma^2 \tau} \right)^{-1} \qquad \text{(S.14)}$$

For non-contiguous intervals ($k > 1$), $R \simeq 0$, and for contiguous intervals ($k = 1$), $R$ remains small.

3

## Multiple timescales

Now we consider the case where rate fluctuations are governed by multiple timescales $\{\tau_i\}$. We define the rate covariance as:

$$\langle \delta\lambda(s)\delta\lambda(s')\rangle = \sum_i \sigma_i^2 \exp\left(-\frac{|s-s'|}{\tau_i}\right) \tag{S.15}$$

Each timescale $\tau_i$ is weighted by a variance $\sigma_i^2$. A given component of the variance $\sigma_i^2$ can be negative, but the total variance $\sigma^2 = \sum_i \sigma_i^2$ is positive.

The spike count variance is then given by:

$$\mathrm{Var}(N) = \mu\Delta + \sum_{i=1}^{n} 2\sigma_i^2\tau_i^2\left[\exp\left(-\frac{\Delta}{\tau_i}\right) - \left(1-\frac{\Delta}{\tau_i}\right)\right] \simeq \Delta\left[\mu + 2\sum_{i\in S}\sigma_i^2\tau_i + \sum_{i\in L}\sigma_i^2\Delta\right] \tag{S.16}$$

where the sets $S$ and $L$ correspond, respectively, to the short and long timescales, namely $S = \{i : |\tau_i| \ll \Delta\}$ and $L = \{i : |\tau_i| \gg \Delta\}$. Unless the absolute variance $|\sigma_i|^2$ of the short timescales is much larger than that of the long timescales, the contribution of the short timescales can be neglected. The spike-count autocovariance is similarly generalized:

$$C(k\Delta) = \sum_i \sigma_i^2\left[2\tau_i\sinh\left(\frac{\Delta}{2\tau_i}\right)\right]^2\exp\left(-\frac{k\Delta}{\tau_i}\right) \simeq \delta_{k1}\sum_{i\in S}\sigma_i^2\tau_i^2 + \sum_{i\in L}\sigma_i^2\Delta^2\exp\left(-\frac{k\Delta}{\tau_i}\right) \tag{S.17}$$

Again, short timescales have negligible contribution to autocovariance unless their variance is much larger than the long timescales.

The autocorrelation is obtained by dividing Equation S.17 by Equation S.16:

$$R(k\Delta) \simeq \frac{\Delta\sum_{i\in L}\sigma_i^2\exp\left(-\frac{k\Delta}{\tau_i}\right)}{\mu + \Delta\sum_{i\in L}\sigma_i^2} \tag{S.18}$$

Note that if time constants are very large, and the value of $k$ is not comparably large, the exponential functions of the corresponding terms are approximately constant. If only one time constant $\tau_j$ shows an appreciable exponential decay in the range of $k$ considered, the autocorrelation is approximately equal to (the sums run over the long timescales only):

$$R(k\Delta) \simeq \left[1 + \frac{\mu}{\sigma_j^2\Delta} + \frac{1}{\Delta}\sum_{i\in S}\tau_i\frac{\sigma_i^2}{\sigma_j^2} + \sum_{i\in L}\frac{\sigma_i^2}{\sigma_j^2}\right]^{-1}\left[\exp\left(-\frac{k\Delta}{\tau_j}\right) + \sum_{i\in L}\frac{\sigma_i^2}{\sigma_j^2}\right] \tag{S.19}$$

4

The last term provides an effective "offset" in the autocorrelation that reflects the contributions of long timescales to spike-rate fluctuations. In line with Equation S.19, we fit the experimental spike-count autocorrelation data with the functional form:

$$R(k\Delta) = A\left[\exp\left(-\frac{k\Delta}{\tau}\right) + B\right] \tag{S.20}$$

These three parameters were used for the fit: $A$ is the amplitude, $\tau$ is the intrinsic timescale, and $B$ is the offset. Equation S.19 supports our interpretation of the autocorrelation offset ($B$ in Equation S.20) as the relative strength of contributions from long timescales compared to that of the moderate intrinsic timescale.

5