

1 Flexible and efficient 2 simulation-based inference for 3 models of decision-making

4 Jan Boelts^{1,2}, Jan-Matthis Lueckmann¹, Richard Gao¹, Jakob H. Macke^{1,3}

*For correspondence:

jan.boelts@uni-tuebingen.de (JB);
jakob.macke@uni-tuebingen.de
(JHM)

5 ¹Machine Learning in Science, Excellence Cluster Machine Learning, University of
6 Tübingen; ²Technical University of Munich; ³Max Planck Institute for Intelligent Systems
7 Tübingen

9 **Abstract** Inferring parameters of computational models that capture experimental data is a
10 central task in cognitive neuroscience. Bayesian statistical inference methods usually require the
11 ability to evaluate the likelihood of the model—however, for many models of interest in cognitive
12 neuroscience, the associated likelihoods cannot be computed efficiently. Simulation-based
13 inference (SBI) offers a solution to this problem by only requiring access to simulations produced
14 by the model. Here, we provide an efficient SBI method for models of decision-making. Our
15 approach, Mixed Neural Likelihood Estimation (MNLE), trains neural density estimators on model
16 simulations to emulate the simulator, and is designed to capture both the continuous (e.g.,
17 reaction times) and discrete (choices) data of decision-making models. The likelihoods of the
18 emulator can then be used to perform Bayesian parameter inference on experimental data using
19 standard approximate inference methods like Markov Chain Monte Carlo sampling. We
20 demonstrate MNLE on two variants of the drift-diffusion model (DDM) and compare its
21 performance to a recently proposed method for SBI on DDMs, called Likelihood Approximation
22 Networks (LANs, Fengler et al. 2021). We show that MNLE is substantially more efficient than
23 LANs: it achieves similar likelihood accuracy with six orders of magnitude fewer training
24 simulations, and is substantially more accurate than LANs when both are trained with the same
25 budget. This enables researchers to train MNLE on custom-tailored models of decision-making,
26 leading to fast iteration of model design for scientific discovery.

28 Introduction

29 Computational modeling is an essential part of the scientific process in cognitive neuroscience:
30 Models are developed from prior knowledge and hypotheses, and compared to experimentally ob-
31 served phenomena (*Churchland and Sejnowski, 1988; McClelland, 2009*). Computational models
32 usually have free parameters which need to be tuned to find those models that capture experi-
33 mental data. This is often approached by searching for single best-fitting parameters using grid
34 search or optimization methods. While this point-wise approach has been used successfully (*Lee*
35 *et al., 2016; Patil et al., 2016*) it can be scientifically more informative to perform Bayesian infer-
36 ence over the model parameters: Bayesian inference takes into account prior knowledge, reveals
37 *all* the parameters consistent with observed data, and thus can be used for quantifying uncer-
38 tainty, hypothesis testing, and model selection (*Lee, 2008; Shiffrin et al., 2008; Lee and Wagen-*
39 *makers, 2014; Schad et al., 2021*). Yet, as the complexity of models used in cognitive neuroscience
40 increases, Bayesian inference becomes challenging for two reasons. First, for many commonly

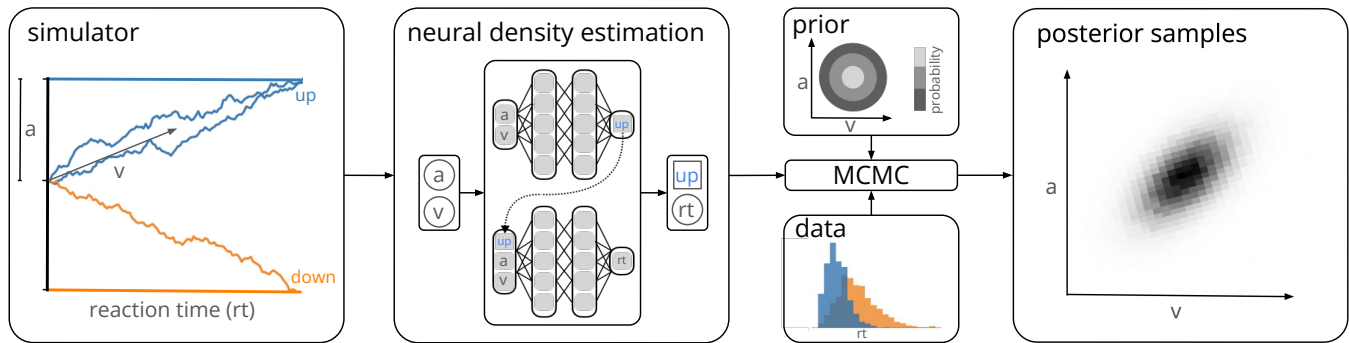


Figure 1. Training a neural density estimator on simulated data to perform parameter inference. Our goal is to perform Bayesian inference on models of decision-making for which likelihoods cannot be evaluated (here a drift-diffusion model for illustration, left). We train a neural density estimation network on synthetic data generated by the model, to provide access to (estimated) likelihoods. Our neural density estimators are designed to account for the mixed data types of decision-making models (e.g. discrete valued choices and continuous valued reaction times, middle). The estimated likelihoods can then be used for inference with standard Markov Chain Monte Carlo (MCMC) methods, i.e., to obtain samples from the posterior over the parameters of the simulator given experimental data (right). Once trained, our method can be applied to flexible inference scenarios like varying number of trials or hierarchical inference without having to retrain the density estimator.

41 used models, computational evaluation of likelihoods is challenging because often no analytical
42 form is available. Numerical approximations of the likelihood are typically computationally expen-
43 sive, rendering standard approximate inference methods like Markov chain Monte Carlo (MCMC)
44 inapplicable. Second, models and experimental paradigms in cognitive neuroscience often induce
45 scenarios in which inference is repeated for varying numbers of experimental trials and changing
46 hierarchical dependencies between model parameters (Lee, 2011). As such, fitting computational
47 models with arbitrary hierarchical structures to experimental data often still requires idiosyncratic
48 and complex inference algorithms.

49 Approximate Bayesian computation (ABC, Sisson et al., 2018) offers a solution to the first chal-
50 lenge by enabling Bayesian inference based on comparing simulated with experimental data, with-
51 out the need to evaluate an explicit likelihood function. Accordingly, various ABC methods have
52 been applied to and developed for models in cognitive neuroscience and related fields (Turner
53 and Van Zandt, 2012, 2018; Palestro et al., 2018a; Kangasrääsiö et al., 2019). However, ABC meth-
54 ods are limited regarding the second challenge because they become inefficient as the number
55 of model parameters increases (Lueckmann et al., 2021) and require generating new simulations
56 whenever the observed data or parameter dependencies change.

57 More recent approaches from the field simulation-based inference (SBI, Cranmer et al., 2020)
58 have the potential to overcome these limitations by using machine learning algorithms such as
59 neural networks. Recently, Fengler et al. (2021) presented an SBI-algorithm for a specific problem
60 in cognitive neuroscience—*inference for drift-diffusion models (DDM)*. They introduced a new ap-
61 proach, called *likelihood approximation networks (LANs)*, which uses neural networks to predict
62 log-likelihoods from data and parameters. The predicted likelihoods can subsequently be used to
63 generate posterior samples using MCMC methods. LANs are trained in a three-step procedure.
64 First, a set of N parameters is generated and for each of the N parameters the model is simulated
65 M times. Second, for each of the N parameters, empirical likelihood targets are estimated from
66 the M model simulations using kernel density estimation (KDE) or empirical histograms. Third,
67 a training dataset consisting of parameters, data points and empirical likelihood targets is con-
68 structed by augmenting the initial set of N parameters by a factor of 1000: for each parameter,
69 1000 data points and empirical likelihood targets are generated from the learned KDE. Finally, su-
70 pervised learning is used to train a neural network to predict log-likelihoods, by minimizing a loss
71 function (the Huber loss) between the network-predicted log-likelihoods and the (log of) the empir-
72 ically estimated likelihoods. Overall, LANs require a large number of model simulations such that
73 the histogram-probability of each possible observed data and for each possible combination of

74 input parameters, can be accurately estimated— $N \cdot M$ model simulations, e.g., $1.5 \times 10^6 \times 10^5$ (150
75 billion) for the examples used in *Fengler et al. (2021)*. The extremely high number of model simula-
76 tions will make it infeasible for most users to run this training themselves, so that there would need
77 to be a repository from which users can download pre-trained LANs. This restricts the application
78 of LANs to a small set of canonical models like drift-diffusion models, and prohibits customization
79 and iteration of models by users. In addition, the high simulation requirement limits this approach
80 to models whose parameters and observations are sufficiently low-dimensional for histograms to
81 be sampled densely.

82 To overcome these limitations, we propose an alternative approach called Mixed Neural Like-
83 lihood Estimation (MNLE). MNLE builds on recent advances in probabilistic machine learning, and
84 in particular on the framework of neural likelihood estimation (*Papamakarios et al., 2019b; Lueck-*
85 *mann et al., 2019*) but is designed to specifically capture the mixed data types arising in models
86 of decision-making, e.g., discrete choices and continuous reaction times. Neural likelihood esti-
87 mation has its origin in classical synthetic likelihood (SL) approaches (*Wood, 2010; Drovandi et al.,*
88 *2018*). Classical SL approaches assume the likelihood of the simulation-based model to be Gaus-
89 sian (so that its moments can be estimated from model simulations) and then use MCMC methods
90 for inference. This approach and various extensions of it have been widely used (*Price et al., 2018;*
91 *Ong et al., 2018; An et al., 2019; Priddle et al., 2021*)—but inherently they need multiple model
92 simulations for each parameter in the MCMC chain to estimate the associated likelihood.

93 *Neural* likelihood approaches instead perform *conditional* density estimation, i.e., they train a
94 neural network to predict the parameters of the approximate likelihood conditioned on the model
95 parameters (e.g., *Papamakarios et al., 2019b; Lueckmann et al., 2019*). By using a conditional den-
96 sity estimator, it is possible to exploit continuity across different model parameters, rather than
97 having to learn a separate density for each individual parameter as in classical SL. Recent advances
98 in conditional density estimation (such as normalizing flows, *Papamakarios et al., 2019a*) further
99 allow lifting the parametric assumptions of classical SL methods and learning flexible conditional
100 density estimators which are able to model a wide range of densities, as well as highly nonlinear de-
101 pendencies on the conditioning variable. In addition, neural likelihood estimators yield estimates
102 of the probability density which are guaranteed to be non-negative and normalized, and which can
103 be both sampled and evaluated, acting as a probabilistic emulator of the simulator (*Lueckmann*
104 *et al., 2019*).

105 Our approach, MNLE, uses neural likelihood estimation to learn an emulator of the simulator.
106 The training phase is a simple two-step procedure: first, a training dataset of N parameters θ
107 is sampled from a proposal distribution and corresponding model simulations \mathbf{x} are generated.
108 Second, the N parameter-data pairs (θ, \mathbf{x}) are directly used to train a conditional neural likelihood
109 estimator to estimate $p(\mathbf{x}|\theta)$. Like for LANs, the proposal distribution for the training data can be *any*
110 distribution over θ , and should be chosen to cover all parameter-values one expects to encounter in
111 empirical data. Thus, the prior distribution used for Bayesian inference constitutes a useful choice,
112 but in principle any distribution that contains the support of the prior can be used. To account
113 for mixed data types, we learn the likelihood estimator as a mixed model composed of one neural
114 density estimator for categorical data and one for continuous data, conditioned on the categorical
115 data. This separation allows us to choose the appropriate neural density estimator for each data
116 type, e.g., a Bernoulli model for the categorical data and a normalizing flow (*Papamakarios et al.,*
117 *2019a*) for the continuous data. The resulting joint density estimator gives access to the likelihood,
118 which enables inference via MCMC methods. See *Figure 1* for an illustration of our approach, and
119 *Methods and Materials* for details.

120 Both LANs and MNLEs allow for flexible inference scenarios common in cognitive neuroscience,
121 e.g., varying number of trials with same underlying experimental conditions or hierarchical infer-
122 ence, and need to be trained only once. However, there is a key difference between the two ap-
123 proaches. LANs use feed-forward neural networks to perform regression from model parameters
124 to empirical likelihood targets obtained from KDE. MNLE instead learns the likelihood directly by

125 performing conditional density estimation on the simulated data without requiring likelihood tar-
126 gets. This makes MNLE by design more simulation efficient than LANs—we demonstrate empir-
127 ically that it can learn likelihood-estimators which are as good or better than those reported in
128 the LAN paper, but using a factor of 1,000,000 fewer simulations (*Fengler et al., 2021*). When us-
129 ing the same simulation-budget for both approaches, MNLE substantially outperforms LAN across
130 several performance metrics. Moreover, MNLE results in a density estimator that is guaranteed
131 to correspond to a valid probability distribution and can also act as an emulator that can gener-
132 ate synthetic data without running the simulator. The simulation-efficiency of MNLEs allows users
133 to explore and iterate on their own models without generating a massive training dataset, rather
134 than restricting their investigations to canonical models for which pre-trained networks have been
135 provided by a central resource. To facilitate this process, we implemented our method as an exten-
136 sion to an open-source toolbox for SBI methods (*Tejero-Cantero et al., 2020*), and provide detailed
137 documentation and tutorials.

138 Results

139 Evaluating the performance of mixed neural likelihood estimation (MNLE) on the 140 drift-diffusion model

141 We first demonstrate the efficiency and performance of MLNEs on a classical model of decision-
142 making, the drift-diffusion model (DDM, *Ratcliff and McKoon, 2008*). The DDM is an influential
143 phenomenological model of a two-alternative perceptual decision-making task. It simulates the
144 evolution of an internal decision variable that integrates sensory evidence until one of two decision
145 boundaries is reached and a choice is made (*Figure 1*, left). The decision variable is modeled with
146 a stochastic differential equation which, in the “simple” DDM version (as used in *Fengler et al.,*
147 *2021*), has four parameters: the drift rate v , boundary separation a , the starting point w of the
148 decision variable, and the non-decision time τ . Given these four parameters $\theta = (v, a, w, \tau)$, a single
149 simulation of the DDM returns data \mathbf{x} containing a choice $c \in \{0, 1\}$ and the corresponding reaction
150 time in seconds $rt \in (\tau, \infty)$: $\mathbf{x} = (c, rt)$.

151 MNLE learns accurate likelihoods with a fraction of the simulation budget

152 The simple version of the DDM is the ideal candidate for comparing the performance of different
153 inference methods because the likelihood of an observation given the parameters, $L(\mathbf{x}|\theta)$, can be
154 calculated analytically (*Navarro and Fuss, 2009*, in contrast to more complicated versions of the
155 DDM, e.g., *Ratcliff and Rouder (1998)*; *Usher and McClelland (2001)*; *Reynolds and Rhodes (2009)*).
156 This enabled us to evaluate MNLE’s performance with respect to the analytical likelihoods and the
157 corresponding inferred posteriors of the DDM, and to compare against that of LANs on a range
158 of simulation-budgets. For MNLE we used a budget of 10^5 simulations (henceforth referred to as
159 MNLE⁵), for LANs we used budgets of 10^5 and 10^8 simulations (LAN⁵, LAN⁸, respectively, trained by
160 us) and the pre-trained version based on 10^{11} simulations (LAN¹¹) provided by *Fengler et al. (2021)*.
161 First, we evaluated the quality of likelihood approximations of MNLE⁵, and compared it to that
162 of LAN^{5,8,11}. Both MNLEs and LANs were in principle able to accurately approximate the likelihoods
163 for both decisions and a wide range of reaction times (see *Figure 2a* for an example, and *Details of*
164 *the numerical comparison*). However, LANs require a much larger simulation budget than MNLE
165 to achieve accurate likelihood approximations, i.e., LANs trained with 10^5 or 10^8 simulations show
166 visible deviations, both in the linear and in log-domain (*Figure 2a*, lines for LAN⁵ and LAN⁸).

167 To quantify the quality of likelihood approximation, we calculated the Huber loss and the mean-
168 squared error (MSE) between the true and approximated likelihoods (*Figure 2b,c*), as well as be-
169 tween the *log*-likelihoods (*Figure 2d,e*). The metrics were calculated as averages over (log-)likelihoods
170 of a fixed observation given 1000 parameters sampled from the prior, repeated for 100 obser-
171 vations simulated from the DDM. For metrics calculated on the untransformed likelihoods (*Fig-*
172 *ure 2b,c*), we found that MNLE⁵ was more accurate than LAN^{5,8,11} on all simulation budgets, show-

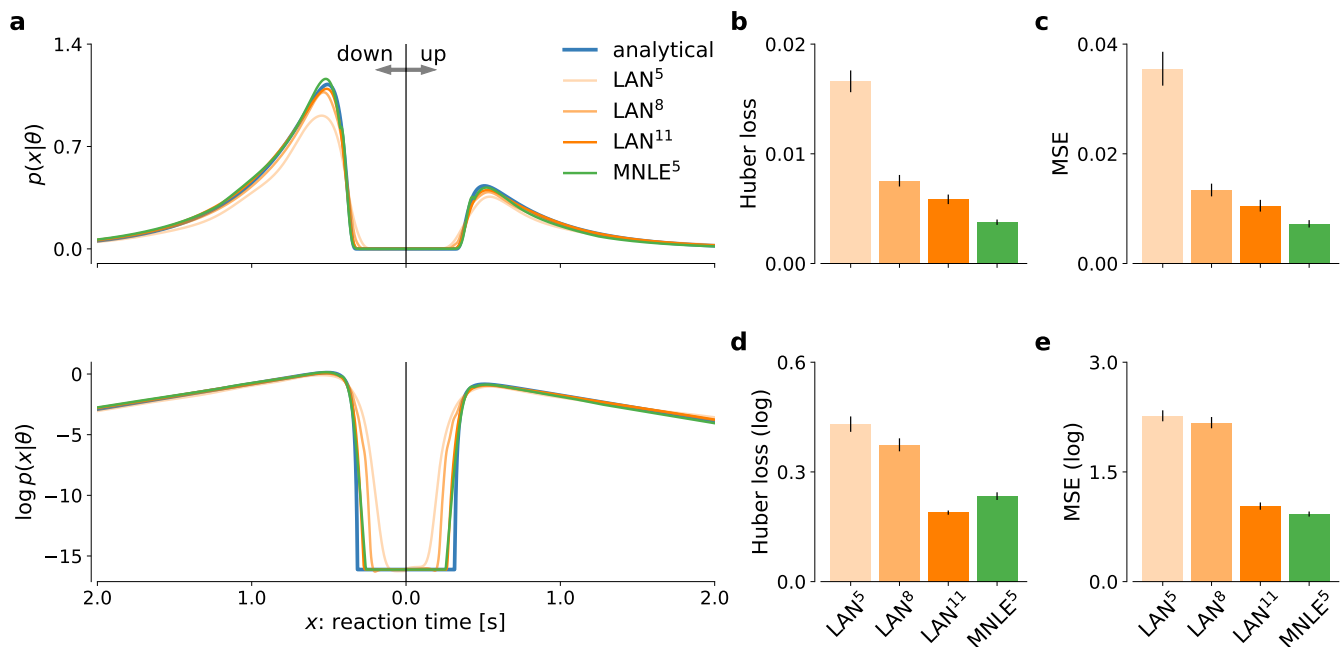


Figure 2. MNLE estimates accurate likelihoods for the drift-diffusion model. The classical drift-diffusion model (DDM) simulates reaction times and choices of a two-alternative decision task and has an analytical likelihood which can be used for comparing the likelihood approximations of MNLE and LAN. We compared MNLE trained with a budget of 10^5 simulations (green, MNLE⁵) to LAN trained with budgets of 10^5 , 10^8 and 10^{11} simulations (shades of orange, LAN^{5,8,11}, respectively). (a) Example likelihood for a fixed parameter θ over a range of reaction times (reaction times for down- and up-choices shown towards the left and right, respectively). Shown on a linear scale (top panel) and a logarithmic scale (bottom panel). (b) Huber loss between analytical and estimated likelihoods calculated for a fixed simulated data point over 1,000 test parameters sampled from the prior, averaged over 100 data points (lower is better). Bar plot error bars show standard error of the mean. (c) Same as in (b), but using mean squared error (MSE) over likelihoods (lower is better). (d) Huber loss on the log-likelihoods (LAN's training loss). (e) MSE on the log-likelihoods.

Figure 2—Figure supplement 1. Examples of synthetic DDM data generated from the MNLE emulator.

173 ing smaller Huber loss than LAN^{5,8,11} in 99, 81 and 66 out of 100 comparisons, and smaller MSE
 174 than LAN^{5,8,11} on 98, 81 and 66 out of 100 comparisons, respectively. Similarly, for the MSE calcu-
 175 lated on the log-likelihoods (Figure 2e), MNLE⁵ achieved smaller MSE than LAN^{5,8,11} on 100, 100
 176 and 75 out of 100 comparisons, respectively. For the Huber loss calculated on the log-likelihoods
 177 (Figure 2d), we found that MNLE⁵ was more accurate than LAN⁵ and LAN⁸, but slightly less accurate
 178 than LAN¹¹, showing smaller Huber loss than LAN^{5,8} in all 100 comparisons, and larger Huber loss
 179 than LAN¹¹ in 62 out of 100 comparisons. All the above pairwise comparisons were significant un-
 180 der the binomial test ($p < 0.01$), but note that these are simulated data and therefore the p-value
 181 can be arbitrarily inflated by increasing the number of comparisons. We also note that the Huber
 182 loss on the log-likelihoods is the loss which is directly optimized by LANs, and thus this comparison
 183 should in theory favor LANs over alternative approaches. Furthermore, the MNLE⁵ results shown
 184 here represent averages over ten random neural network initializations (five of which achieved
 185 smaller Huber loss than LAN¹¹), whereas the LAN¹¹ results are based on a single pre-trained net-
 186 work. Finally, we also investigated MNLE's property to act as an emulator of the simulator and
 187 found the synthetic reaction times and choices generated from the MNLE emulator to match cor-
 188 responding data simulated from the DDM accurately (see Figure 2—Figure Supplement 1 and Ap-
 189 pendix 1).

190 When using the learned likelihood estimators for inference with MCMC methods, their evalu-
 191 ation speed can also be important because MCMC often requires thousands of likelihood evalu-
 192 ations. We found that evaluating MNLE for a batch of 100 trials and ten model parameters (as
 193 used during MCMC) took 4.14 ± 0.04 ms (mean over 100 repetitions \pm standard error of the mean),

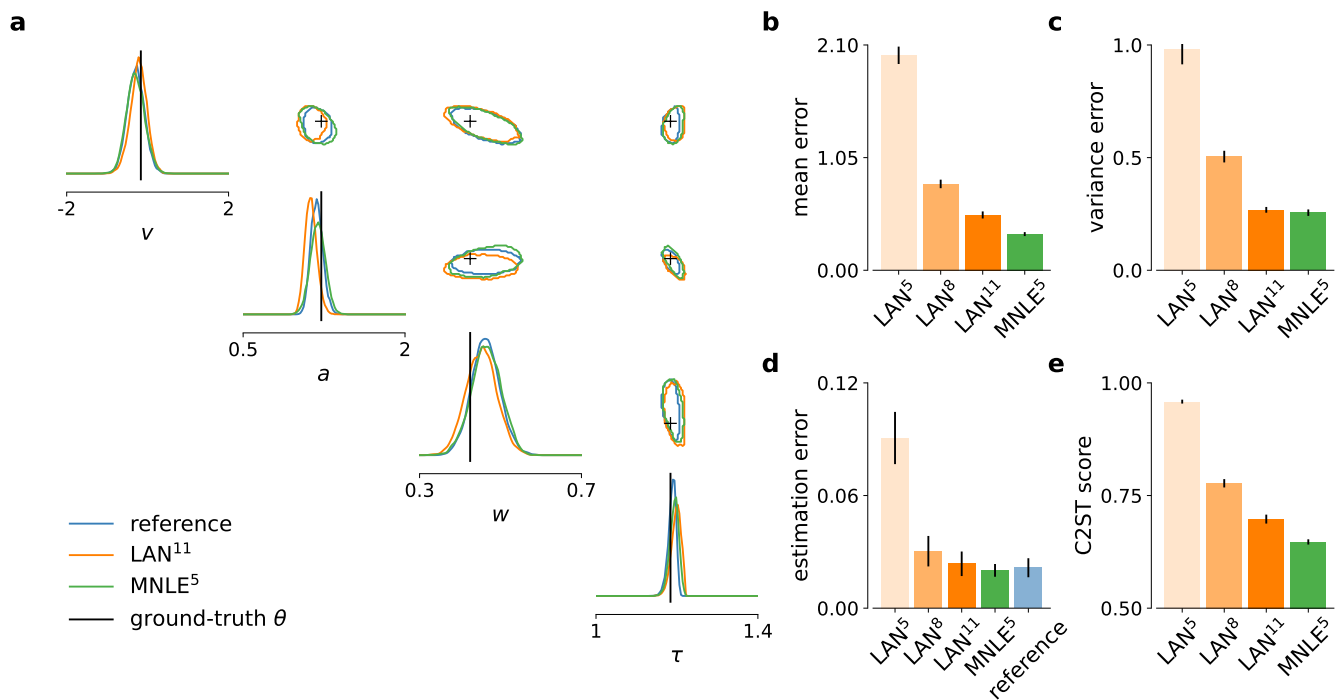


Figure 3. MNLE infers accurate posteriors for the drift-diffusion model. Posteriors were obtained given 100-trial i.i.d. observations with MCMC using analytical (i.e., reference) likelihoods, or those approximated using LAN^{5,8,11} trained with simulation budgets 10^{5,8,11}, respectively, and MNLE⁵ trained with a budget of 10⁵ simulations. (a) Posteriors given an example observation generated from the prior and the simulator, shown as 95% contour lines in a corner-plot, i.e., one-dimension marginal (diagonal) and all pairwise two-dimensional marginals (upper triangle). (b) Difference in posterior sample mean of approximate (LAN^{5,8,11}, MNLE⁵) and reference posteriors (normalized by reference posterior standard deviation, lower is better). (c) Same as in (b) but for posterior sample variance (normalized by reference posterior variance, lower is better). (d) Parameter estimation error measured as mean squared error (MSE) between posterior sample mean and the true underlying parameters (smallest possible error is given by reference posterior performance shown in blue). (e) Classification 2-sample test (C2ST) score between approximate (LAN^{5,8,11}, MNLE⁵) and reference posterior samples (0.5 is best). All bar plots show metrics calculated from 100 repetitions with different observations; error bars show standard error of the mean.

Figure 3–Figure supplement 1. Inference accuracy metrics for individual model parameters.

Figure 3–Figure supplement 2. Example posteriors and parameter recovery for LAN⁵ and LAN⁸.

Figure 3–Figure supplement 3. Inference accuracy metrics for different numbers of observed trials.

194 compared to 1.02 ± 0.03 ms for LANs, i.e., MNLE incurred a larger computational foot-print at evalu-
 195 ation time. Note that these timings are based on an improved implementation of LANs compared
 196 to the one originally presented in *Fengler et al. (2021)*, and evaluation times can depend on the
 197 implementation, compute infrastructure and parameter settings (see *Details of the numerical com-*
 198 *parison and Discussion*). In summary, we found that MNLE trained with 10⁵ simulations performed
 199 substantially better than LANs trained with 10⁵ or 10⁸ simulations, and similarly well or better than
 200 LANs trained with 10¹¹ simulations, on all likelihood approximation accuracy metrics.

201 MNLE enables accurate flexible posterior inference with MCMC

202 In the previous section we showed that MNLE requires substantially fewer training simulations than
 203 LANs to approximate the likelihood accurately. To investigate whether these likelihood-estimates
 204 were accurate enough to support accurate parameter inference, we evaluated the quality of the
 205 resulting posteriors, using a framework for benchmarking SBI algorithms (*Lueckmann et al., 2021*).
 206 We used the analytical likelihoods of the simple DDM to obtain reference posteriors for 100 differ-
 207 ent observations, via MCMC sampling. Each observation consisted of 100 independent and identi-
 208 cally distributed (i.i.d.) trials simulated with parameters sampled from the prior (see *Figure 3a* for
 209 an example, details in *Methods and Materials*). We then performed inference using MCMC based

210 on the approximate likelihoods obtained with MNLE (10^5 budget, MNLE⁵) and the ones obtained
211 with LAN for each of the three simulation budgets (LAN^{5,8,11}, respectively).

212 Overall, we found that the likelihood approximation performances presented above were re-
213 flected in the inference performances: MNLE⁵ performed substantially better than LAN⁵ and LAN⁸,
214 and equally well or better than LAN¹¹ (**Figure 3b-d**). In particular, MNLE⁵ approximated the poste-
215 rior mean more accurately than LAN^{5,8,11} (**Figure 3b**), being more accurate than LAN^{5,8,11} in 100,
216 90, and 67 out of 100 comparisons, respectively. In terms of posterior variance, MNLE⁵ performed
217 better than LAN^{5,8} and on par with LAN¹¹ (**Figure 3c**), being more accurate than LAN^{5,8,11} in 100,
218 93, ($p \ll 0.01$, binomial test) and 58 ($p = 0.13$) out of 100 pairwise comparisons, respectively.

219 Additionally, we measured the parameter estimation accuracy as the mean squared error be-
220 tween the posterior mean and the ground-truth parameters underlying the observed data. We
221 found that MNLE⁵ estimation error was indistinguishable from that of the reference posterior, and
222 that LAN performance was similar only for the substantially larger simulation budget of LAN¹¹ (**Fig-
223 ure 3c**), with MNLE being closer to reference performance than LAN^{5,8,11} in 100, 91, and 66 out
224 of 100 comparisons, respectively (all $p < 0.01$). Note that all three metrics were reported as av-
225 erages over the four parameter dimensions of the DDM to keep the visualizations compact, and
226 that this average did not affect the results qualitatively. We report metrics for each dimension in
227 **Figure 3—Figure Supplement 1**, as well as additional inference accuracy results for smaller LAN
228 simulation budgets (**Figure 3—Figure Supplement 2**) and for different numbers of observed trials
229 (**Figure 3—Figure Supplement 3**).

230 Finally, we used the classifier 2-sample test (C2ST, *Lopez-Paz and Oquab, 2017; Lueckmann
231 et al., 2021*) to quantify the similarity between the estimated and reference posterior distributions.
232 The C2ST is defined to be the error-rate of a classification algorithm which aims to classify whether
233 samples belong to the true or the estimated posterior. Thus, it ranges from 0.5 (no difference
234 between the distributions, the classifier is at chance level), to 1.0 (the classifier can perfectly distin-
235 guish the two distributions). We note that the C2ST is a highly sensitive measure of discrepancy
236 between two multivariate-distributions—e.g. if the two distributions differ in *any* dimension, the
237 C2ST will be close to 1 even if all other dimensions match perfectly. We found that neither of the
238 two approaches was able to achieve perfect approximations, but that MNLE⁵ achieved lower C2ST
239 scores compared to LAN^{5,8,11} on all simulation budgets (**Figure 3e**): mean C2ST score LAN^{5,8,11},
240 0.96, 0.78, 0.70 vs. MNLE⁵, 0.65, with MNLE⁵ showing a better score than LAN^{5,8,11} on 100, 91, and
241 68 out of 100 pairwise comparisons, respectively (all $p < 0.01$). In summary, MNLE achieves more
242 accurate recovery of posterior means than LANs, similar or better recovery of posterior variances,
243 and overall more accurate posteriors (as quantified by C2ST).

244 **MNLE posteriors have uncertainties which are well-calibrated**

245 For practical applications of inference, it is often desirable to know how well an inference proce-
246 dure can recover the ground-truth parameters, and whether the uncertainty-estimates are well-
247 calibrated, (*Cook et al., 2006*), i.e., that the *uncertainty* estimates of the posterior are balanced, and
248 neither over-confident nor under-confident. For the DDM, we found that the posteriors inferred
249 with MNLE and LANs (when using LAN¹¹) recovered the ground-truth parameters accurately (in
250 terms of posterior means, **Figure 3d** and **Figure 4a**)—in fact, parameter recovery was similarly ac-
251 curate to using the ‘true’ analytical likelihoods, indicating that much of the residual error is due to
252 stochasticity of the observations, and not the inaccuracy of the likelihood approximations.

253 To assess posterior calibration, we used simulation-based calibration (SBC, *Talts et al., 2018*).
254 The basic idea of SBC is the following: If one repeats the inference with simulations from many dif-
255 ferent prior samples, then, with a well-calibrated inference method, the combined samples from
256 all the inferred posteriors should be distributed according to the prior. One way to test this is
257 to calculate the rank of each ground-truth parameter (samples from the prior) under its corre-
258 sponding posterior, and to check whether all the ranks follow a uniform distribution. SBC results
259 indicated that MNLE posteriors were as well-calibrated as the reference posteriors, i.e., the empir-

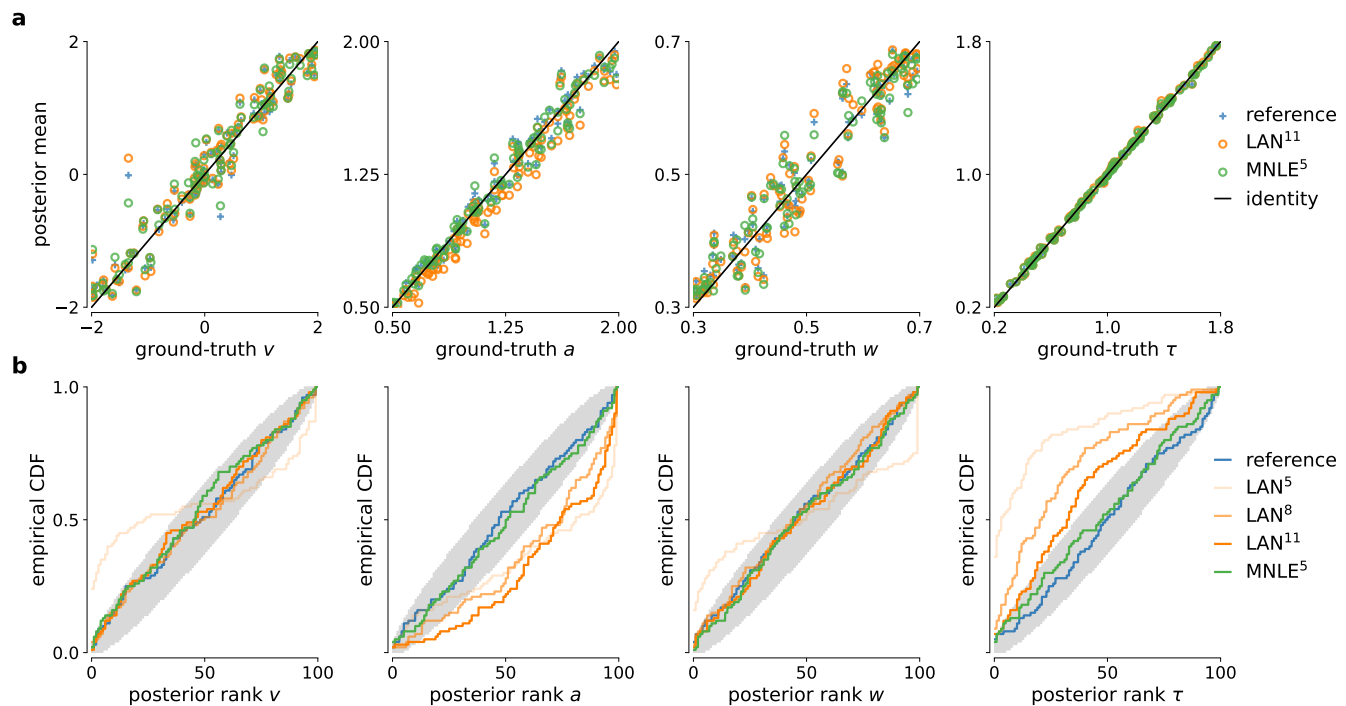


Figure 4. Parameter recovery and posterior uncertainty calibration for the DDM. (a) Underlying ground-truth DDM parameters plotted against the sample mean of posterior samples inferred with the analytical likelihoods (reference, blue crosses), LAN (orange circles) and MNLE (green circles), for 100 different observations. Markers close to diagonal indicate good recovery of ground-truth parameters; circles on top of blue reference crosses indicate accurate posterior means. (b) Simulation-based calibration results showing empirical cumulative density functions (CDF) of the ground-truth parameters ranked under the inferred posteriors calculated from 100 different observations. A well-calibrated posterior must have uniformly distributed ranks, as indicated by the area shaded gray. Shown for reference posteriors (blue), LAN posteriors obtained with increasing simulation budgets (shades of orange, LAN^{5,8,11}) and MNLE posterior (green, MNLE⁵), and for each parameters separately (v , a , w and τ).

Figure 4—Figure supplement 1. Parameter recovery for different numbers of observed trials.

Figure 4—Figure supplement 2. Simulation-based calibration results for different numbers of observed trials.

260 ical cumulative density functions (CDF) of the ranks were close to that of a uniform distribution
 261 (Figure 4b)—thus, on this example, MNLE inferences would likely be of similar quality compared to
 262 using the analytical likelihoods. When trained with the large simulation budget of 10^{11} simulations,
 263 LANs, too appeared to recover most of the ground-truth parameters well. However, SBC detected
 264 a systematic underestimation of the parameter a and overestimation of the parameter τ , and this
 265 bias increased for the smaller simulation budgets of LAN⁵ and LAN⁸ (Figure 4b, see the deviation
 266 below and above the desired uniform distribution of ranks, respectively).

267 The results so far (i.e., Figure 3, Figure 4) indicate that both LAN¹¹ and MNLE⁵ lead to similar pa-
 268 rameter recovery, but only MNLE⁵ leads to posteriors which were well-calibrated for all parameters.
 269 These results were obtained using a scenario with 100 i.i.d. trials. When increasing the number of
 270 trials (e.g., to 1000 trials), posteriors become very concentrated around the ground-truth value. In
 271 that case, while the posteriors overall identified the ground-truth parameter value very well (Fig-
 272 ure 4—Figure Supplement 1c), even small deviations in the posteriors can have large effects on the
 273 posterior metrics (Figure 3—Figure Supplement 3). This effect was also detected by SBC, showing
 274 systematic biases for some parameters (Figure 4—Figure Supplement 2c). For MNLE, we found
 275 that these biases were smaller, and furthermore that it was possible to mitigate this effect by infer-
 276 ring the posterior using ensembles, e.g., by combining samples inferred with five MNLEs trained
 277 with identical settings but different random initialization (see Appendix 1 for details). These results
 278 show the utility of using SBC as a tool to test posterior coverage, especially τ when studying models

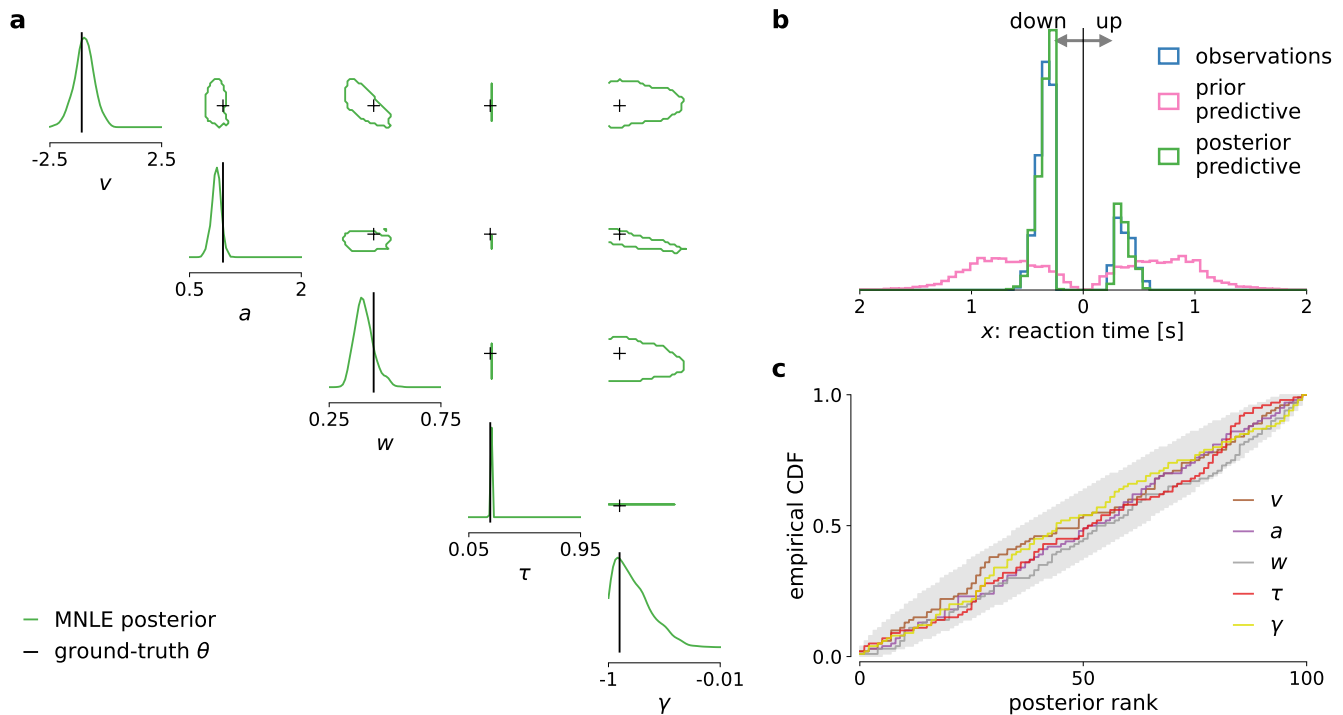


Figure 5. MNLE infers accurate posteriors for the DDM with collapsing bounds. Posterior samples were obtained given 100-trial observations simulated from the DDM with linearly collapsing bounds, using MNLE and MCMC. **(a)** Approximate posteriors shown as 95% contour lines in a corner-plot of one-dimensional (diagonal) and two-dimensional (upper triangle) marginals, for a representative 100-trial observation simulated from the DDM. **(b)** Reaction times and choices simulated from the ground-truth parameters (blue) compared to those simulated given parameters sampled from the prior (prior predictive distribution, purple) and from the MNLE posterior shown in (a) (posterior predictive distribution, green). **(c)** Simulation-based calibration results showing empirical cumulative density functions (CDF) of the ground-truth parameters ranked under the inferred posteriors, calculated from 100 different observations. A well-calibrated posterior must have uniformly distributed ranks, as indicated by the area shaded gray. Shown for each parameters separately (v , a , w , τ and γ).

279 for which reference posteriors are not available, as we demonstrate in the next section.

280 MNLE infers well-calibrated, predictive posteriors for a DDM with collapsing bounds

281 MNLE was designed to be applicable to models for which evaluation of the likelihood is not practical so that standard inference tools cannot be used. To demonstrate this, we applied MNLE to a
 282 variant of the DDM for which analytical likelihoods are not available (note, however, that numerical
 283 approximation of likelihoods for this model would be possible, see e.g., *Shinn et al., 2020*, and
 284 *Methods and Materials* for details). This DDM variant simulates a decision variable like the simple
 285 DDM used above, but with linearly collapsing instead of constant decision boundaries (see e.g.,
 286 *Hawkins et al., 2015; Palestro et al., 2018b*). The collapsing bounds are incorporated with an additional
 287 parameter γ indicating the slope of the decision boundary, such that $\theta = (a, v, w, \tau, \gamma)$ (see
 288 *Details of the numerical comparison*).

290 We tested inference with MNLE on the DDM with linearly collapsing bound using observations
 291 comprised of 100 i.i.d. trials simulated with parameters sampled from the prior. Using the same
 292 MNLE training and MCMC settings as above, we found that posterior density concentrated around
 293 the underlying ground-truth parameters (see *Figure 5a*), suggesting that MNLE learned the under-
 294 lying likelihood accurately. To assess inference quality systematically without needing reference
 295 posteriors, we performed posterior predictive checks by running simulations with the inferred pos-
 296 teriors samples and comparing them to observed (simulated) data, and checked posterior calibra-
 297 tion properties using SBC. We found that the inferred posteriors have good predictive performance,
 298 i.e., data simulated from the inferred posterior samples accurately matched the observed data (*Fig-*

299 *ure 5b*), and that their uncertainties are well-calibrated as quantified by the SBC results (*Figure 5c*).
300 Overall, this indicated that MNLE accurately inferred the posterior of this intractable variant of the
301 DDM.

302 Discussion

303 Statistical inference for computational models in cognitive neuroscience can be challenging be-
304 cause models often do not have tractable likelihood functions. The recently proposed LAN method
305 (*Fengler et al., 2021*) performs SBI for a subset of such models (DDMs) by training neural networks
306 with model simulations to approximate the intractable likelihood. However, LANs require large
307 amounts of training data, restricting its usage to canonical models. We proposed an alternative ap-
308 proached called mixed neural likelihood estimation (MNLE), a synthetic neural likelihood method
309 which is tailored to the data-types encountered in many models of decision-making.

310 Our comparison on a tractable example problem used in *Fengler et al. (2021)* showed that
311 MNLE performed on par with LANs using six orders of magnitude fewer model simulations for
312 training. While *Fengler et al. (2021)* discuss that LANs were not optimized for simulation efficiency
313 and that it might be possible to reduce the required model simulations, we emphasize that the
314 difference in simulation-efficiency is due to an inherent property of LANs. For each parameter in
315 the training data, LANs require empirical likelihood targets that have to be estimated by building
316 histograms or kernel density estimates from thousands of simulations. MNLE, instead, performs
317 conditional density estimation without the need of likelihood targets and can work with only one
318 simulation per parameter. Because of these conceptual differences, we expect the substantial
319 performance advantage of MNLE to be robust to the specifics of the implementation.

320 After the networks are trained, the time needed for each evaluation determines the speed of in-
321 ference. In that respect, both LANs and MNLEs are conceptually similar in that they require a single
322 forward-pass through a neural network for each evaluation, and we found MNLE and the original
323 implementation of LANs to require comparable computation times. However, evaluation time will
324 depend, e.g., on the exact network architecture, software framework and computing infrastructure
325 used. Code for a new `PyTorch` implementation of LANs has recently been released and improved
326 upon the evaluation speed original implementation we compared to. While this new implementa-
327 tion made LAN significantly faster for a single forward-pass, we observed that the resulting time
328 difference with the MCMC-settings used here was only on the order of minutes, whereas the differ-
329 ence in simulation time for LAN¹¹ vs MNLE⁵ was on the order of days. The exact timings will always
330 be implementation specific and whether or not these differences are important will depend on
331 the application at hand. In a situation where iteration over model design is required (i.e., custom
332 DDMs), an increase in training efficiency on the order of days would be advantageous.

333 There exist a number of approaches with corresponding software packages for estimating pa-
334 rameters of cognitive neuroscience models, and DDMs in particular. However, these approaches
335 either only estimate single best-fitting parameters (*Voss and Voss, 2007; Wagenmakers et al., 2007;*
336 *Chandrasekaran and Hawkins, 2019; Heathcote et al., 2019; Shinn et al., 2020*) or, if they perform
337 fully Bayesian inference, can only produce Gaussian approximations to posteriors (*Feltgen and*
338 *Daunizeau, 2021*), or are restricted to variants of the DDM for which the likelihood can be evalu-
339 ated (*Wiecki et al., 2013*, HDDM). A recent extension of the HDDM toolbox includes LANs, thereby
340 combining HDDM's flexibility with LAN's ability to perform inference without access to the likeli-
341 hood function (but this remains restricted to variants of the DDM for which LAN can be pre-trained).
342 In contrast, MNLE can be applied to any simulation-based model with intractable likelihoods and
343 mixed data type-outputs. Here, we focused on the direct comparison to LANs based on variants
344 of the DDM. We note that these models have a rather low-dimensional observation structure (as
345 common in many cognitive neuroscience models), and that our examples did not include additional
346 parameter structure, e.g., stimulus encoding parameters, which would increase the dimensionality
347 of the learning problem. However, other variants of neural density estimation have been applied
348 successfully to a variety of problems with higher dimensionality (see e.g. *Gonçalves et al., 2020;*

349 *Lueckmann et al., 2021; Glöckler et al., 2021; Dax et al., 2022*). Therefore, we expect MNLE to be
350 applicable to other simulation-based problems with higher-dimensional observation structure and
351 parameter spaces, and to scale more favourably than LANs. Like for any neural network-based ap-
352 proach, applying MNLE to different inference problems may require selecting different architecture
353 and training hyperparameters settings, which may induce additional computational training costs.
354 To help with this process, we adopted default settings which have been shown to work well on a
355 large range of SBI benchmarking problems (*Lueckmann et al., 2021*), and we integrated MNLE into
356 the established `sbi` python package that provides well-documented implementations for training-
357 and inference performance of SBI algorithms.

358 Several extensions to classical synthetic likelihood (SL) approaches have addressed the problem
359 of a bias in the likelihood approximation due to the strong parametric assumptions, i.e., Gaussian-
360 ity, the use of summary statistics, or finite-sample biases (*Price et al., 2018; Ong et al., 2018; van*
361 *Opheusden et al., 2020*). MNLE builds on flexible neural likelihood estimators, e.g., normalizing
362 flows, and does not require summary statistics for a low-dimensional simulator like the DDM, so
363 would be less susceptible to these first two biases. It could be subject to biases resulting from
364 the estimation of the log-likelihoods from a finite number of simulations. In our numerical experi-
365 ments, and for the simulation-budgets we used, we did not observe biased inference results. We
366 speculate that the ability of neural density estimators to pool data across multiple parameter set-
367 tings (rather than using only data from a specific parameter set, like in classical synthetic likelihood
368 methods) mitigates finite-sample effects.

369 MNLE is a SBI method which uses neural density estimators to estimate likelihoods. Alterna-
370 tives to neural likelihood estimation include neural posterior estimation (NPE, *Papamakarios and*
371 *Murray, 2016; Lueckmann et al., 2017; Greenberg et al., 2019*, which uses conditional density es-
372 timation to learn the posterior directly) and neural ratio estimation (NRE, *Hermans et al., 2020;*
373 *Durkan et al., 2020*, which uses classification to approximate the likelihood-to-evidence ratio to
374 then perform MCMC). These approaches could in principle be applied here as well, however, they
375 are not as well suited for the flexible inference scenarios common in decision-making models as
376 MNLE: NPE by design does not allow for flexible inference scenarios but needs to be retrained
377 because the posterior changes with changing number of trials or changing hierarchical inference
378 setting; and NRE, performing ratio- and not density estimation, would not provide an emulator of
379 the simulator.

380 Regarding future research directions, MNLE has the potential to become more simulation effi-
381 cient by using weight sharing between the discrete and the continuous neural density estimators
382 (rather than to use separate neural networks, as we did here). Moreover, for high-dimensional
383 inference problems in which slice sampling-based MCMC might struggle, MNLE could be used in
384 conjunction with gradient-based MCMC methods like Hamiltonian Monte Carlo (HMC, *Neal et al.,*
385 *2011; Hoffman et al., 2014*), or variational inference as recently proposed by *Wiqvist et al. (2021)*
386 and *Glöckler et al. (2021)*. With MNLE's full integration into the `sbi` package, both gradient-based
387 MCMC methods from Pyro (*Bingham et al., 2019*), and variational inference for SBI (SNVI, *Glöckler*
388 *et al., 2021*) are available as inference methods for MNLE (a comparison of HMC and SNVI to slice
389 sampling-MCMC on one example observation resulted in indistinguishable posterior samples). Fi-
390 nally, using its emulator property (see Appendix 1), MNLE could be applied in an active learning
391 setting for highly expensive simulators in which new simulations are chosen adaptively accord-
392 ing to a acquisition function in a Bayesian optimization framework (*Gutmann and Corander, 2016;*
393 *Lueckmann et al., 2019; Järvenpää et al., 2019*).

394 In summary, MNLE enables flexible and efficient inference of parameters of models in cognitive
395 neuroscience with intractable likelihoods. The training efficiency and flexibility of the neural density
396 estimators used overcome the limitations of LANs (*Fengler et al., 2021*). Critically, these features
397 enable researchers to develop customized models of decision-making and not just apply existing
398 models to new data. We implemented our approach as an extension to a public `sbi` python package
399 with detailed documentation and examples to make it accessible for practitioners.

400 Methods and Materials

401 Mixed neural likelihood estimation

402 Mixed neural likelihood estimation (MNLE) extends the framework of neural likelihood estimation
403 (*Papamakarios et al., 2019b; Lueckmann et al., 2019*) to be applicable to simulation-based models
404 with mixed data types. It learns a parametric model $q_{\psi}(\mathbf{x}|\theta)$ of the intractable likelihood $p(\mathbf{x}|\theta)$
405 defined implicitly by the simulation-based model. The parameters ψ are learned with training data
406 $\{\theta_n, \mathbf{x}_n\}_{1:N}$ comprised of model parameters θ_n and their corresponding data simulated from the
407 model $\mathbf{x}_n|\theta_n \sim p(\mathbf{x}|\theta_n)$. The parameters are sampled from a proposal distribution over parameters
408 $\theta_n \sim p(\theta)$. The proposal distribution could be any distribution, but it determines the parameter
409 regions for which the density estimator will be good in estimating likelihoods. Thus, the prior,
410 or a distribution that contains the support of the prior (or even all priors which one expects to
411 use in the future) constitutes a useful choice. After training, the emulator can be used to generate
412 synthetic data $\mathbf{x}|\theta \sim q_{\psi}(\mathbf{x}|\theta)$ given parameters, and to evaluate the synthetic likelihood $q_{\psi}(\mathbf{x}|\theta)$ given
413 experimentally observed data. Finally, the synthetic likelihood can be used to obtain posterior
414 samples via

$$p(\theta|\mathbf{x}) \propto q_{\psi}(\mathbf{x}|\theta)p(\theta), \quad (1)$$

415 through approximate inference with MCMC. Importantly, the training is amortized, i.e., the em-
416 ulator can be applied to new experimental data without retraining (running MCMC is still required).

417 We tailored MNLE to simulation-based models that return mixed data, e.g., in form of reaction
418 times rt and (usually categorical) choices c as for the DDM. Conditional density estimation with nor-
419 malizing flows has been proposed for continuous random variables (*Papamakarios et al., 2019a*),
420 or discrete random variables (*Tran et al., 2019*), but not for mixed data. Our solution for estimat-
421 ing the likelihood of mixed data is to simply factorize the likelihood into continuous and discrete
422 variables,

$$p(rt, c|\theta) = p(rt|\theta, c) p(c|\theta), \quad (2)$$

423 and use two separate neural likelihood estimators: A discrete one q_{ψ_c} to estimate $p(c|\theta)$ and a
424 continuous one $q_{\psi_{rt}}$ to estimate $p(rt|\theta, c)$. We defined q_{ψ_c} to be a Bernoulli model and use a neural
425 network to learn the Bernoulli probability ρ given parameters θ . For $q_{\psi_{rt}}$ we used a conditional
426 neural spline flow (*Durkan et al., 2019*) to learn the density of rt given a parameter θ and choice
427 c . The two estimators are trained separately using the same training data (see Neural network
428 architecture, training and hyperparameters for details). After training, the full neural likelihood
429 can be constructed by multiplying the likelihood estimates q_{ψ_c} and $q_{\psi_{rt}}$ back together:

$$q_{\psi_c, \psi_{rt}}(rt, c|\theta) = q_{\psi_c}(c|\theta) q_{\psi_{rt}}(rt|c, \theta). \quad (3)$$

430 Note that, as the second estimator $q_{\psi_{rt}}(rt|c, \theta)$ is conditioned on the choice c , our likelihood-
431 model can account for statistical dependencies between choices and reaction times. The neural
432 likelihood can then be used to approximate the intractable likelihood defined by the simulator,
433 e.g., for inference with MCMC. Additionally, it can be used to sample synthetic data given model
434 parameters, without running the simulator (see The emulator property of MNLE).

435 Relation to LAN

436 Neural likelihood estimation can be much more simulation efficient than previous approaches be-
437 cause it exploits continuity across the parameters by making the density estimation conditional.
438 *Fengler et al. (2021)*, too, aim to exploit continuity across the parameter space by training a neural
439 network to predict the value of the likelihood function from parameters θ and data \mathbf{x} . However,
440 the difference to neural likelihood estimation is that they do not use the neural network for density
441 estimation directly, but instead do classical neural network-based regression on likelihood targets.

442 Crucially, the likelihood targets first have to be obtained for each parameter in the training data set.
443 To do so, one has to perform density estimation using KDE (as proposed by *Turner et al., 2015*)
444 or empirical histograms for every parameter separately. Once trained, LANs do indeed exploit
445 the continuity across the parameter space (they can predict log-likelihoods given unseen data and
446 parameters), however, they do so at a very high simulation cost: For a training data set of N param-
447 eters, they perform N times KDE based on M simulations each¹, resulting in an overall simulation
448 budget of $N \cdot M$ ($N = 1.5$ million and $M = 100,000$ for “pointwise” LAN approach).

449 **Details of the numerical comparison**

450 The comparison between MNLE and LAN is based on the drift-diffusion model (DDM). The DDM
451 simulates a decision variable X as a stochastic differential equation with parameters $\theta = (v, a, w, \tau)$:

$$452 \quad dX_{t+\tau} = vdt + dW, \quad X_\tau = w, \quad (4)$$

452 where W a Wiener noise process. The priors over the parameters are defined to be uniform: $v \sim$
453 $\mathcal{U}(-2, 2)$ is the drift, $a \sim \mathcal{U}(0.5, 2)$ the boundary separation, $w \sim \mathcal{U}(0.3, 0.7)$ the initial offset, $\tau \sim$
454 $\mathcal{U}(0.2, 1.8)$ the non-decision time. A single simulation from the model returns a choice $c \in \{0, 1\}$
455 and the corresponding reaction time in seconds $rt \in (\tau, \infty)$.

456 For this version of the DDM the likelihood of an observation (c, rt) given parameters θ , $L(c, rt|\theta)$,
457 can be calculated analytically (*Navarro and Fuss, 2009*). To simulate the DDM and calculate ana-
458 lytical likelihoods we used the approach and the implementation proposed by *Drugowitsch (2016)*.
459 We numerically confirmed that the simulations and the analytical likelihoods match those obtained
460 from the research code provided by *Fengler et al. (2021)*.

461 To run LANs, we downloaded the neural network weights of the pre-trained models from the
462 repository mentioned in *Fengler et al. (2021)*. The budget of training simulations used for the LANs
463 was 1.5×10^{11} (1.5 million training data points, each obtained from KDE based on 10^5 simulations).
464 We only considered the approach based on training a multilayer-perceptron (MLP) on single-trial
465 likelihoods (“pointwise approach”, *Fengler et al., 2021*). At a later stage of our study we performed
466 additional experiments to evaluate the performance of LANs trained at smaller simulation budgets,
467 e.g., for 10^5 and 10^8 simulations. For this analysis we used an updated implementation of LANs
468 based on PyTorch that was provided by the authors. We used the training routines and default
469 settings provided with that implementation. To train LANs at the smaller budgets we used the
470 following splits of budgets into number of parameter settings drawn from the prior, and number
471 of simulations per parameter setting used for fitting the KDE: for the 10^5 budget we used 10^2 and 10^3 ,
472 respectively (we ran experiments splitting the other way around, but results were slightly better
473 for this split); for the 10^8 budget we used an equal split of 10^4 each (all code publicly available, see
474 Code availability).

475 To run MNLE, we extended the implementation of neural likelihood estimation in the *sbi* tool-
476 box (*Tejero-Cantero et al., 2020*). All comparisons were performed on a single AMD Ryzen Thread-
477 ripper 1920X 12-Core processor with 2.2GHz and the code is publicly available (see Code availabil-
478 ity).

479 For the DDM variant with linearly collapsing decision boundaries, the boundaries were parametrized
480 by the initial boundary separation, a , and one additional parameter, γ , indicating the slope with
481 which the boundary approaches zero. This resulted in a five-dimensional parameter space for
482 which we used the same prior as above, plus the an additional uniform prior for the slope $\gamma \sim$
483 $\mathcal{U}(-1.0, 0)$. To simulate this DDM variant, we again used the Julia package by *Drugowitsch (2016)*,
484 but we note that for this variant no analytical likelihoods are available. While it would be possi-
485 ble to approximate the likelihoods numerically using the Fokker-Planck equations (see, e.g., *Shinn*
486 *et al., 2020*), this would usually involve a trade-off between computation time and accuracy as well
487 as design of bespoke solutions for individual models, and was not pursued here.

¹For models with categorical output, i.e., all decision-making models, KDE is performed separately for each choice.

488 **Flexible Bayesian inference with MCMC**

489 Once the MNLE is trained, it can be used for MCMC to obtain posterior samples $\theta \sim p(\theta|\mathbf{x})$ given
490 experimentally observed data \mathbf{x} . To sample from posteriors via MCMC, we transformed the param-
491 eters to unconstrained space, used slice sampling (Neal, 2003), and initialized ten parallel chains
492 using sequential importance sampling (Papamakarios et al., 2019b), all as implemented in the `sbi`
493 toolbox. We ran MCMC with identical settings for MNLE and LAN.

494 Importantly, performing MNLE and then using MCMC to obtain posterior samples allows for
495 flexible inference scenarios because the form of \mathbf{x} is not fixed. For example, when the model pro-
496 duces trial-based data that satisfy the i.i.d. assumption, e.g., a set of reaction times and choices
497 $\mathbf{X} = \{rt, c\}_{i=1}^N$ in a drift-diffusion model, then MNLE allows to perform inference given varying num-
498 bers of trials, without retraining. This is achieved by training MNLE based on single-trial likelihoods
499 once and then combining multiple trials into the joint likelihood only when running MCMC:

$$p(\theta|\mathbf{X}) \propto \prod_{i=1}^N q(rt_i, c_i|\theta) p(\theta). \quad (5)$$

500 Similarly, one can use the neural likelihood to perform hierarchical inference with MCMC, all with-
501 out the need for retraining (see Hermans et al., 2020; Fengler et al., 2021, for examples).

502 Stimulus- and inter-trial dependencies

503 Simulation-based models in cognitive neuroscience often depend not only on a set of parameters
504 θ , but additionally on (a set of) stimulus variables s which are typically given as part of the exper-
505 imental conditions. MNLE can be readily adapted to this scenario by generating simulated data
506 with multiple stimulus variables, and including them as additional inputs to the network during in-
507 ference. Similarly, MNLE could be adapted to scenarios in which the i.i.d. assumption across trials
508 as used above (see Eq. Flexible Bayesian inference with MCMC) does not hold. Again, this would be
509 achieved by adapting the model-simulator accordingly. For example, when inferring parameters
510 θ of a DDM for which the outcome of the current trial i additionally depends on current stimulus
511 variables s_i as well as on previous stimuli s_j and responses r_j , then one would implement the DDM
512 simulator as a function $f(\theta; s_{i-T}, \dots, s_i; r_{i-T}, \dots, r_{i-1})$ (where T is a history parameter) and then learn
513 the underlying likelihood by emulating f with MNLE.

514 **Neural network architecture, training and hyperparameters**

515 Architecture

516 For the architecture of the Bernoulli model we chose a feed-forward neural network that takes
517 parameters θ as input and predicts the Bernoulli probability ρ of the corresponding choices. For
518 the normalizing flow we used the neural spline flow architecture (NSF, Durkan et al., 2019). NSFs
519 use a standard normal base distribution and transform it using several modules of monotonic
520 rational-quadratic splines whose parameters are learned by invertible neural networks. Using an
521 unbounded base distribution for modeling data with bounded support, e.g., reaction time data $rt \in$
522 $(0, \infty)$, can be challenging. To account for this, we tested two approaches: We either transformed
523 the reaction time data to logarithmic space to obtain an unbounded support ($\log rt \in (-\infty, \infty)$), or
524 we used a log-normal base distribution with rectified (instead of linear) tails for the splines (see
525 Durkan et al., 2019, for details and Architecture and training hyperparameters for the architecture
526 settings used)

527 Training

528 The neural network parameters ψ_c and ψ_{rt} were trained using the maximum likelihood loss and
529 the Adam optimizer (Kingma and Ba, 2015). As proposal distribution for the training dataset we
530 used the prior over DDM parameters. Given a training data set of parameters, choices and reaction
531 times $\{\theta_i, (c_i, rt_i)\}_{i=1}^N$ with $\theta_i \sim p(\theta)$; $(c_i, rt_i) \sim \text{DDM}(\theta_i)$, we minimized the negative log-probability of

532 the model. In particular, using the same training data, we trained the Bernoulli choice model by
533 minimizing

$$-\frac{1}{N} \sum_{i=1}^N \log q_{\psi_c}(c_i | \theta_i), \quad (6)$$

534 and the neural spline flow by minimizing

$$-\frac{1}{N} \sum_{i=1}^N \log q_{\psi_{rt}}(rt | c_i, \theta_i). \quad (7)$$

535 Training was performed with code and training hyperparameter settings provided in the `sbi` tool-
536 box (*Tejero-Cantero et al., 2020*).

537 Hyperparameters

538 MNLE requires a number of hyperparameter choices regarding the neural network architectures,
539 e.g., number of hidden layers, number of hidden units, number of stacked NSF transforms, kind
540 of base distribution, among others (*Durkan et al., 2019*). With our implementation building on the
541 `sbi` package we based our hyperparameter choices on the default settings provided there. This
542 resulted in likelihood accuracy similar to LAN, but longer evaluation times due to the complexity
543 of the underlying normalizing flow architecture.

544 To reduce evaluation time of MNLE, we further adapted the architecture to the example model
545 (DDM). In particular, we ran a cross-validation of the hyperparameters relevant for evaluation time,
546 i.e., number of hidden layers, hidden units, NSF transforms, spline bins, and selected those that
547 were optimal in terms of Huber loss and mean-squared error between the approximate and the
548 *analytical* likelihoods, as well as evaluation time. This resulted in an architecture with performance
549 *and* evaluation time similar to LANs (more details in Appendix *Architecture and training hyperpa-*
550 *rameters*). The cross-validation relied on access to the analytical likelihoods which is usually not
551 given in practice, e.g., for simulators with intractable likelihoods. However, we note that in cases
552 without access to analytical likelihoods a similar cross-validation can be performed using quality
553 measures other than the difference to the analytical likelihood, e.g., by comparing the observed
554 data with synthetic data and synthetic likelihoods provided by MNLE.

555 Acknowledgments

556 We thank Luigi Acerbi, Michael Deistler, Alexander Fengler, Michael Frank and Ingeborg Wenger for
557 discussions and comments on a preliminary version of the manuscript. We also acknowledge and
558 thank the Python (*Van Rossum and Drake Jr, 1995*) and Julia (*Bezanson et al., 2017*) communities
559 for developing the tools enabling this work, including `DifferentialEquations.jl` (*Rackauckas and*
560 *Nie, 2017*), `DiffModels.jl` (*Drugowitsch, 2016*), `NumPy` (*Harris et al., 2020*), `pandas` (*pandas develop-*
561 *ment team, 2020*), `Pyro` (*Bingham et al., 2019*), `PyTorch` (*Paszke et al., 2019*), `sbi` (*Tejero-Cantero*
562 *et al., 2020*), `sbibm` (*Lueckmann et al., 2021*) and `Scikit-learn` (*Pedregosa et al., 2011*).

563 References

- 564 **An Z**, South LF, Nott DJ, Drovandi CC. Accelerating Bayesian synthetic likelihood with the graphical lasso. *Journal*
565 *of Computational and Graphical Statistics*. 2019; 28(2):471–475.
- 566 **Bezanson J**, Edelman A, Karpinski S, Shah VB. Julia: A fresh approach to numerical computing. *SIAM review*.
567 2017; 59(1):65–98.
- 568 **Bingham E**, Chen JP, Jankowiak M, Obermeyer F, Pradhan N, Karaletsos T, Singh R, Szerlip P, Horsfall P, Good-
569 man ND. Pyro: Deep Universal Probabilistic Programming. *Journal of Machine Learning Research*. 2019;
570 20(1):973–978.
- 571 **Chandrasekaran C**, Hawkins GE. ChaRT: An R toolbox for modeling choices and response times in decision-
572 making tasks. *Journal of neuroscience methods*. 2019; 328:108432.

- 573 **Churchland PS**, Sejnowski TJ. Perspectives on cognitive neuroscience. *Science*. 1988; 242(4879):741–745.
- 574 **Cook SR**, Gelman A, Rubin DB. Validation of software for Bayesian models using posterior quantiles. *Journal*
575 *of Computational and Graphical Statistics*. 2006; 15(3):675–692.
- 576 **Cranmer K**, Brehmer J, Louppe G. The frontier of simulation-based inference. *Proceedings of the National*
577 *Academy of Sciences*. 2020; doi: [10.1073/pnas.1912789117](https://doi.org/10.1073/pnas.1912789117).
- 578 **Dax M**, Green SR, Gair J, Deistler M, Schölkopf B, Macke JH. Group equivariant neural posterior estimation. In:
579 *International Conference on Learning Representations*; 2022. <https://openreview.net/forum?id=u6s8dSporO8>.
- 580 **Drovandi CC**, Grazian C, Mengersen K, Robert C. Approximating the Likelihood in ABC. *Handbook of approxi-*
581 *mate bayesian computation*. 2018; p. 321–368.
- 582 **Drugowitsch J**. Fast and accurate Monte Carlo sampling of first-passage times from Wiener diffusion models.
583 *Scientific reports*. 2016; 6(1):1–13.
- 584 **Durkan C**, Bekasov A, Murray I, Papamakarios G. Neural spline flows. *Advances in Neural Information Process-*
585 *ing Systems*. 2019; 32:7511–7522.
- 586 **Durkan C**, Murray I, Papamakarios G. On contrastive learning for likelihood-free inference. In: *International*
587 *Conference on Machine Learning* PMLR; 2020. p. 2771–2781.
- 588 **Feltgen Q**, Daunizeau J. An overcomplete approach to fitting drift-diffusion decision models to trial-by-trial
589 data. *Frontiers in artificial intelligence*. 2021; 4:23.
- 590 **Fengler A**, Govindarajan LN, Chen T, Frank MJ. Likelihood approximation networks (LANs) for fast inference of
591 simulation models in cognitive neuroscience. *eLife*. 2021; 10:e65074.
- 592 **Glöckler M**, Deistler M, Macke JH. Variational methods for simulation-based inference. In: *International Con-*
593 *ference on Learning Representations*; 2021. .
- 594 **Gonçalves PJ**, Lueckmann JM, Deistler M, Nonnenmacher M, Öcal K, Bassetto G, Chintaluri C, Podlaski WF, Had-
595 dad SA, Vogels TP, Greenberg DS, Macke JH. Training deep neural density estimators to identify mechanistic
596 models of neural dynamics. *eLife*. 2020; doi: [10.7554/eLife.56261](https://doi.org/10.7554/eLife.56261).
- 597 **Greenberg D**, Nonnenmacher M, Macke J. Automatic Posterior Transformation for Likelihood-Free Inference.
598 In: *Proceedings of the 36th International Conference on Machine Learning*, vol. 97 of *Proceedings of Machine*
599 *Learning Research* PMLR; 2019. p. 2404–2414.
- 600 **Gutmann MU**, Corander J. Bayesian optimization for likelihood-free inference of simulator-based statistical
601 models. *The Journal of Machine Learning Research*. 2016; 17(1):4256–4302.
- 602 **Harris CR**, Millman KJ, van der Walt SJ, Gommers R, Virtanen P, Cournapeau D, Wieser E, Taylor J, Berg S, Smith
603 NJ, et al. Array programming with NumPy. *Nature*. 2020; 585(7825):357–362.
- 604 **Hawkins GE**, Forstmann BU, Wagenmakers EJ, Ratcliff R, Brown SD. Revisiting the evidence for collapsing
605 boundaries and urgency signals in perceptual decision-making. *Journal of Neuroscience*. 2015; 35(6):2476–
606 2484.
- 607 **Heathcote A**, Lin YS, Reynolds A, Strickland L, Gretton M, Matzke D. Dynamic models of choice. *Behavior*
608 *research methods*. 2019; 51(2):961–985.
- 609 **Hermans J**, Begy V, Louppe G. Likelihood-free MCMC with Approximate Likelihood Ratios. In: *Proceedings of*
610 *the 37th International Conference on Machine Learning*, vol. 98 of *Proceedings of Machine Learning Research*
611 *PMLR*; 2020. .
- 612 **Hoffman MD**, Gelman A, et al. The No-U-Turn sampler: adaptively setting path lengths in Hamiltonian Monte
613 Carlo. *J Mach Learn Res*. 2014; 15(1):1593–1623.
- 614 **Järvenpää M**, Gutmann MU, Pleska A, Vehtari A, Marttinen P, et al. Efficient acquisition rules for model-based
615 approximate Bayesian computation. *Bayesian Analysis*. 2019; 14(2):595–622.
- 616 **Kangasrääsiö A**, Jokinen JP, Oulasvirta A, Howes A, Kaski S. Parameter inference for computational cognitive
617 models with Approximate Bayesian Computation. *Cognitive science*. 2019; 43(6):e12738.
- 618 **Kingma DP**, Ba J. Adam: A Method for Stochastic Optimization. In: *Proceedings of the 3rd International Conference*
619 *on Learning Representations, ICLR*; 2015. .

- 620 **Lee HS**, Betts S, Anderson JR. Learning problem-solving rules as search through a hypothesis space. *Cognitive*
621 *science*. 2016; 40(5):1036–1079.
- 622 **Lee MD**. Three case studies in the Bayesian analysis of cognitive models. *Psychonomic Bulletin & Review*. 2008;
623 15(1):1–15.
- 624 **Lee MD**. How cognitive modeling can benefit from hierarchical Bayesian models. *Journal of Mathematical*
625 *Psychology*. 2011; 55(1):1–7.
- 626 **Lee MD**, Wagenmakers EJ. *Bayesian cognitive modeling: A practical course*. Cambridge university press; 2014.
- 627 **Lopez-Paz D**, Oquab M. Revisiting Classifier Two-Sample Tests. In: *5th International Conference on Learning*
628 *Representations, ICLR*; 2017. .
- 629 **Lueckmann JM**, Bassetto G, Karaletsos T, Macke JH. Likelihood-free inference with emulator networks. In:
630 *Proceedings of The 1st Symposium on Advances in Approximate Bayesian Inference*, vol. 96 of Proceedings of
631 *Machine Learning Research PMLR*; 2019. p. 32–53.
- 632 **Lueckmann JM**, Boelts J, Greenberg D, Goncalves P, Macke J. Benchmarking Simulation-Based Inference. In:
633 Banerjee A, Fukumizu K, editors. *Proceedings of The 24th International Conference on Artificial Intelligence and*
634 *Statistics*, vol. 130 of Proceedings of Machine Learning Research PMLR; 2021. p. 343–351.
- 635 **Lueckmann JM**, Goncalves PJ, Bassetto G, Öcal K, Nonnenmacher M, Macke JH. Flexible statistical inference
636 for mechanistic models of neural dynamics. *Advances in Neural Information Processing Systems*. 2017; 30.
- 637 **McClelland JL**. The place of modeling in cognitive science. *Topics in Cognitive Science*. 2009; 1(1):11–38.
- 638 **Navarro DJ**, Fuss IG. Fast and accurate calculations for first-passage times in Wiener diffusion models. *Journal*
639 *of mathematical psychology*. 2009; 53(4):222–230.
- 640 **Neal RM**. Slice sampling. *Annals of Statistics*. 2003; p. 705–741.
- 641 **Neal RM**, et al. MCMC using Hamiltonian dynamics. *Handbook of markov chain monte carlo*. 2011; 2(11):2.
- 642 **Ong VM**, Nott DJ, Tran MN, Sisson SA, Drovandi CC. Variational Bayes with synthetic likelihood. *Statistics and*
643 *Computing*. 2018; 28(4):971–988.
- 644 **van Opheusden B**, Acerbi L, Ma WJ. Unbiased and efficient log-likelihood estimation with inverse binomial
645 sampling. *PLoS computational biology*. 2020; 16(12):e1008483.
- 646 **Palestro JJ**, Sederberg PB, Osth AF, Van Zandt T, Turner BM. *Likelihood-free methods for cognitive science*.
647 Springer; 2018.
- 648 **Palestro JJ**, Weichart E, Sederberg PB, Turner BM. Some task demands induce collapsing bounds: Evidence
649 from a behavioral analysis. *Psychonomic bulletin & review*. 2018; 25(4):1225–1248.
- 650 **Papamakarios G**, Murray I. Fast ϵ -free Inference of Simulation Models with Bayesian Conditional Density
651 Estimation. In: *Advances in Neural Information Processing Systems 29* Curran Associates, Inc.; 2016.p. 1028–
652 1036.
- 653 **Papamakarios G**, Nalisnick E, Rezende DJ, Mohamed S, Lakshminarayanan B. Normalizing flows for probabilis-
654 tic modeling and inference. *arXiv preprint arXiv:191202762*. 2019; .
- 655 **Papamakarios G**, Sterratt D, Murray I. Sequential Neural Likelihood: Fast Likelihood-free Inference with Au-
656 toregressive Flows. In: *Proceedings of the 22nd International Conference on Artificial Intelligence and Statistics*
657 *(AISTATS)*, vol. 89 of Proceedings of Machine Learning Research PMLR; 2019. p. 837–848.
- 658 **Paszke A**, Gross S, Massa F, Lerer A, Bradbury J, Chanan G, Killeen T, Lin Z, Gimelshein N, Antiga L, Desmaison
659 A, Kopf A, Yang E, DeVito Z, Raison M, Tejani A, Chilamkurthy S, Steiner B, Fang L, Bai J, et al. PyTorch: An
660 Imperative Style, High-Performance Deep Learning Library. In: *Advances in Neural Information Processing*
661 *Systems 32* Curran Associates, Inc.; 2019. p. 8024–8035.
- 662 **Patil U**, Hanne S, Burchert F, De Bleser R, Vasisht S. A computational evaluation of sentence processing deficits
663 in aphasia. *Cognitive Science*. 2016; 40(1):5–50.
- 664 **Pedregosa F**, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, Blondel M, Prettenhofer P, Weiss R,
665 Dubourg V, Vanderplas J, Passos A, Cournapeau D, Brucher M, Perrot M, Duchesnay E. Scikit-learn: Machine
666 Learning in Python. *Journal of Machine Learning Research*. 2011; 12:2825–2830.

- 667 **Price LF**, Drovandi CC, Lee A, Nott DJ. Bayesian synthetic likelihood. *Journal of Computational and Graphical*
668 *Statistics*. 2018; 27(1):1–11.
- 669 **Priddle JW**, Sisson SA, Frazier DT, Turner I, Drovandi C. Efficient Bayesian synthetic likelihood with whitening
670 transformations. *Journal of Computational and Graphical Statistics*. 2021; p. 1–14.
- 671 **Rackauckas C**, Nie Q. DifferentialEquations.jl – A Performant and Feature-Rich Ecosystem for Solving Differen-
672 tial Equations in Julia. *The Journal of Open Research Software*. 2017; 5(1). doi: [10.5334/jors.151](https://doi.org/10.5334/jors.151).
- 673 **Ratcliff R**, McKoon G. The diffusion decision model: theory and data for two-choice decision tasks. *Neural*
674 *computation*. 2008; 20(4):873–922.
- 675 **Ratcliff R**, Rouder JN. Modeling response times for two-choice decisions. *Psychological science*. 1998; 9(5):347–
676 356.
- 677 **Reynolds AM**, Rhodes CJ. The Lévy flight paradigm: random search patterns and mechanisms. *Ecology*. 2009;
678 90(4):877–887.
- 679 **Schad DJ**, Betancourt M, Vasishth S. Toward a principled Bayesian workflow in cognitive science. *Psychological*
680 *methods*. 2021; 26(1):103.
- 681 **Shiffrin RM**, Lee MD, Kim W, Wagenmakers EJ. A survey of model evaluation approaches with a tutorial on
682 hierarchical Bayesian methods. *Cognitive Science*. 2008; 32(8):1248–1284.
- 683 **Shinn M**, Lam NH, Murray JD. A flexible framework for simulating and fitting generalized drift-diffusion models.
684 *ELife*. 2020; 9:e56938.
- 685 **Sisson SA**, Y F, A BM. Overview of ABC. In: *Handbook of Approximate Bayesian Computation* CRC Press, Taylor &
686 Francis Group; 2018.
- 687 **Talts S**, Betancourt M, Simpson D, Vehtari A, Gelman A. Validating Bayesian inference algorithms with
688 simulation-based calibration. arXiv preprint arXiv:180406788. 2018; .
- 689 **pandas development team T**, pandas-dev/pandas: Pandas. Zenodo; 2020. [https://doi.org/10.5281/zenodo.](https://doi.org/10.5281/zenodo.3509134)
690 [3509134](https://doi.org/10.5281/zenodo.3509134), doi: [10.5281/zenodo.3509134](https://doi.org/10.5281/zenodo.3509134).
- 691 **Tejero-Cantero A**, Boelts J, Deistler M, Lueckmann JM, Durkan C, Gonçalves PJ, Greenberg DS, Macke JH. sbi: A
692 toolkit for simulation-based inference. *Journal of Open Source Software*. 2020; 5(52):2505. [https://doi.org/](https://doi.org/10.21105/joss.02505)
693 [10.21105/joss.02505](https://doi.org/10.21105/joss.02505), doi: [10.21105/joss.02505](https://doi.org/10.21105/joss.02505).
- 694 **Tran D**, Vafa K, Agrawal K, Dinh L, Poole B. Discrete flows: Invertible generative models of discrete data. *Ad-*
695 *vances in Neural Information Processing Systems*. 2019; 32:14719–14728.
- 696 **Turner BM**, Van Maanen L, Forstmann BU. Informing cognitive abstractions through neuroimaging: the neural
697 drift diffusion model. *Psychological review*. 2015; 122(2):312.
- 698 **Turner BM**, Van Zandt T. A tutorial on approximate Bayesian computation. *Journal of Mathematical Psychology*.
699 2012; 56(2):69–85.
- 700 **Turner BM**, Van Zandt T. Approximating Bayesian inference through model simulation. *Trends in cognitive*
701 *sciences*. 2018; 22(9):826–840.
- 702 **Usher M**, McClelland JL. The time course of perceptual choice: the leaky, competing accumulator model. *Psy-*
703 *chological review*. 2001; 108(3):550.
- 704 **Van Rossum G**, Drake Jr FL. Python tutorial. Centrum voor Wiskunde en Informatica Amsterdam, The Nether-
705 lands; 1995.
- 706 **Voss A**, Voss J. Fast-dm: A free program for efficient diffusion model analysis. *Behavior research methods*.
707 2007; 39(4):767–775.
- 708 **Wagenmakers EJ**, Van Der Maas HL, Grasman RP. An EZ-diffusion model for response time and accuracy.
709 *Psychonomic bulletin & review*. 2007; 14(1):3–22.
- 710 **Wiecki TV**, Sofer I, Frank MJ. HDDM: Hierarchical Bayesian estimation of the drift-diffusion model in Python.
711 *Frontiers in neuroinformatics*. 2013; 7:14.

- 712 **Wiqvist S**, Frellsen J, Picchini U. Sequential Neural Posterior and Likelihood Approximation. arXiv preprint
713 arXiv:210206522. 2021; .
- 714 **Wood SN**. Statistical inference for noisy nonlinear ecological dynamic systems. *Nature*. 2010; 466(7310):1102–
715 1104.

716 **Appendix 1**

717 **Code availability**

718 We implemented MNLE as part of the open source package for SBI, `sbi`, available at <https://github.com/mackelab/sbi>. Code for reproducing the results presented here, and tutorials on how to apply
719 MNLE to other simulators using `sbi` can be found at <https://github.com/mackelab/mnle-for-ddms>.

721 **Architecture and training hyperparameters**

722 For the Bernoulli neural network we used three hidden layers with ten units each and sigmoid acti-
723 vation functions. For the neural spline flow architecture (*Durkan et al., 2019*) we transformed the
724 reaction time data to the log-domain, used a standard normal base distribution, two spline trans-
725 forms with five bins each and conditioning networks with three hidden layers and ten hidden units
726 each, and rectified linear unit activation functions. The neural network training was performed us-
727 ing the `sbi` package with the following settings: learning rate 0.0005; training batch size 100; 10% of
728 training data as validation data, stop training after 20 epochs without validation loss improvement.

729 **The emulator property of MNLE**

730 Being based on the neural likelihood estimation framework, MNLE naturally returns an emulator
731 of the simulator that can be sampled to generate synthetic data without running the simulator.
732 We found that the synthetic data generated by MNLE accurately matched the data we obtained by
733 running the DDM simulator (*Figure 2—Figure Supplement 1*). This has several potential benefits: it
734 can help with evaluating the performance of the density estimator, it enables almost instantaneous
735 data generation (one forward pass in the neural network) even if the simulator is computationally
736 expensive, and it gives full access to the internals of the emulator, e.g., to gradients w.r.t. to data
737 or parameters.

738 There is variant of the LAN approach which allows for sampling synthetic data as well: In the
739 “Histogram-approach” (*Fengler et al., 2021*) LANs are trained with a convolutional neural network
740 (CNN) architecture using likelihood targets in form of two-dimensional empirical histograms. The
741 output of the CNN is a probability distribution over a discretized version of the data space which
742 can, in principle, be sampled to generate synthetic DDM choices and reaction times. However, the
743 accuracy of this emulator property of CNN-LANs is limited by the number of bins used to approxi-
744 mate the continuous data space (e.g., 512 bins for the examples shown in *Fengler et al. (2021)*).

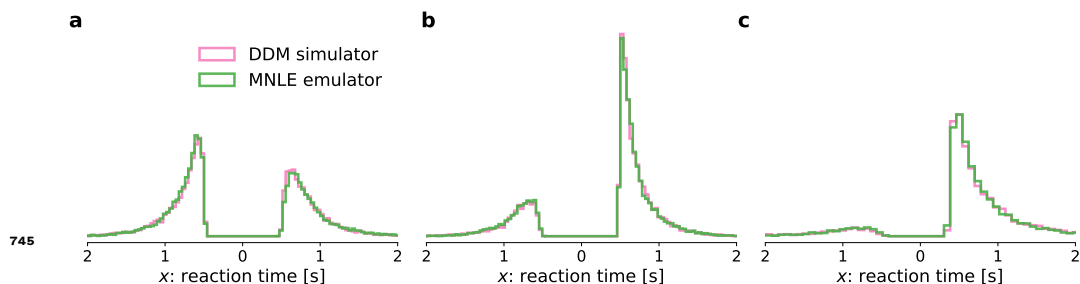


Figure 2-Figure supplement 1. Comparison of simulated DDM data and synthetic data sampled from the MNLE emulator. Histograms of reaction times from 1000 i.i.d. trials generated from three different parameters sampled from the prior (panel **a**, **b**, **c**) using the original DDM simulator (purple) and the emulator obtained from MNLE (green). “Down” choices are shown to the left of zero and “up” choices to the right of zero.

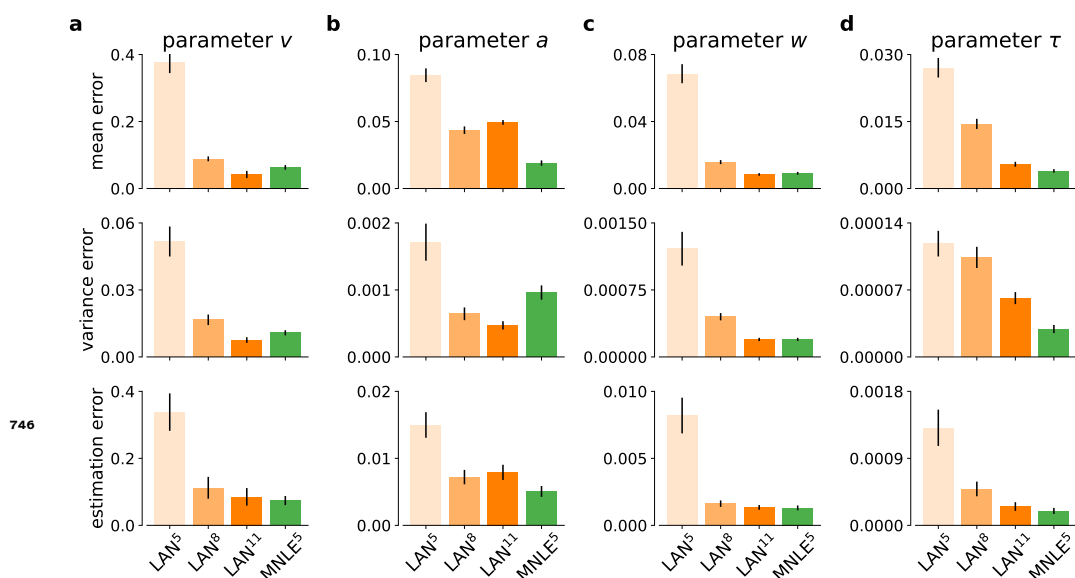


Figure 3-Figure supplement 1. DDM inference accuracy metrics for individual model parameters. Inference accuracy given a 100-trial observation, measured as posterior mean accuracy (first row of panels), posterior variance (second row) and parameter estimation error (third row), shown in absolute terms, for each of the four DDM parameters separately (panels **a**, **b**, **c**, and **d**, respectively), and for each simulation budgets of LAN^(5,8,11) (shades of orange) and for MNLE⁵ trained with 10⁵ simulations (green). Bars show the mean metric over 100 different observations, error bars show standard error of the mean.

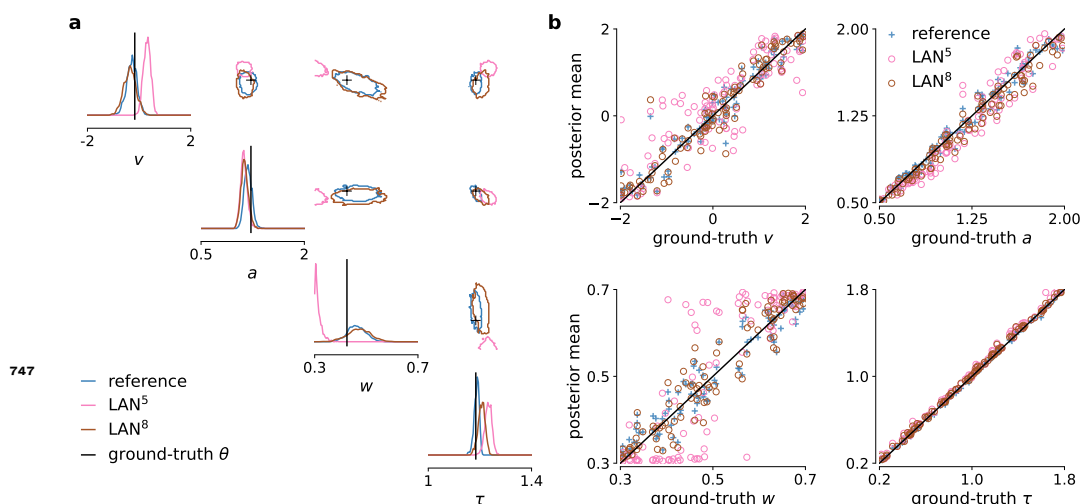


Figure 3-Figure supplement 2. DDM example posteriors and parameter recovery for LANs trained with smaller simulation budgets. (a) Posterior samples given 100-trial example observation, obtained with MCMC using LAN approximate likelihoods trained based on 10^5 (LAN⁵) and 10^8 (LAN⁸) simulations (LAN⁸), and with the analytical likelihoods (reference). (b) Parameter recovery of LAN and the reference posterior shown as posterior sample means against the underlying ground-truth parameters.

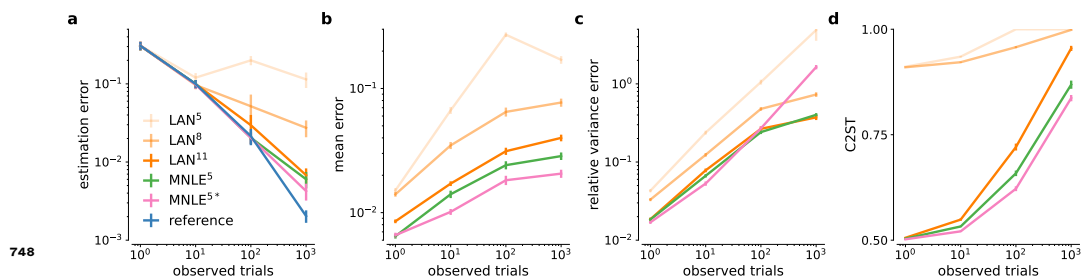


Figure 3-Figure supplement 3. DDM inference accuracy metrics for different numbers of observed trials. Parameter estimation error (a), absolute posterior mean error (b), relative posterior variance error (c) and C2ST scores (d) shown for LAN with increasing simulation budgets (shades of orange, LAN^{5,8,11}), MNLE trained with 10^5 simulations (green), and MNLE ensembles (purple). Metrics were calculated from 10,000 posterior samples and with respect to the reference posterior, for 100 different observations. Error bars show standard error of the mean.

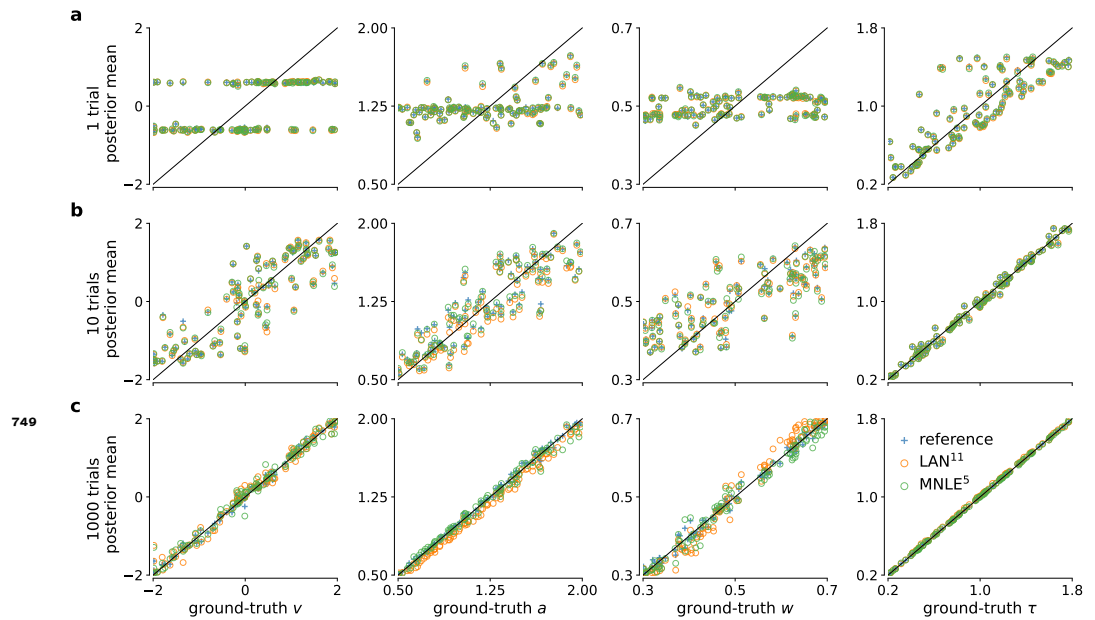


Figure 4-Figure supplement 1. DDM parameter recovery for different number of observed trials. True underlying DDM parameters plotted against posterior sample means for 1, 10, and 1000 of observed i.i.d. trial(s) (in rows) and for the four DDM parameters v , a , w and τ (in columns). Calculated from 10,000 posterior samples obtained with MCMC using the reference (blue), LAN¹¹ (orange) and the MNLE⁵ (green) likelihoods. Black line shows the identity function indicating perfect recovery.

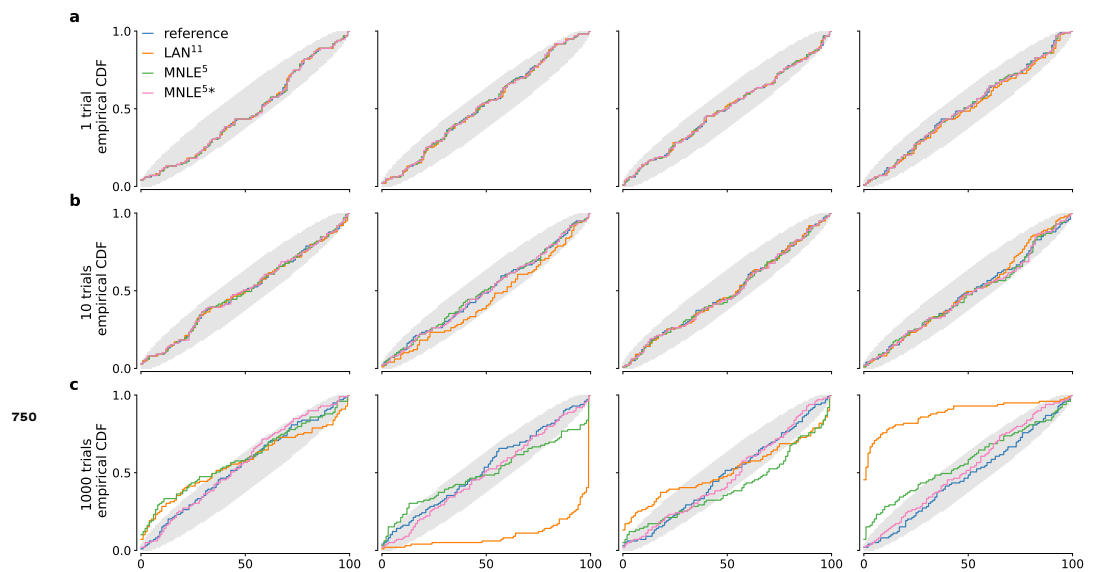


Figure 4-Figure supplement 2. DDM simulation-based calibration results for different numbers of observed trials. SBC results in form empirical conditional density functions of the ranks of ground-truth parameters under the approximate posterior samples. We compared posterior samples based on analytical likelihoods (blue), LAN¹¹ (orange), MNLE⁵ (green), and an ensemble of five MLNEs (MNLE^{5*}, purple); for each of the four parameters of the DDM and for 1, 10, 1000 observed trials (panel **a**, **b**, and **c**, respectively). Grey areas show random deviations expected under uniformly distributed ranks (ideal case).