

Analog Diffusion Models

Jiaqi Chu

Jiaqi.Chu@microsoft.com

Microsoft Research <https://orcid.org/0009-0008-4744-334X>

Heiner Kremer

Microsoft Research

Fabian Falck

Microsoft Research

Grace Brennan

Microsoft research <https://orcid.org/0000-0001-7081-0509>

Burcu Canakci

Microsoft Research

James Clegg

Microsoft Research Ltd <https://orcid.org/0009-0002-3428-7161>

Daniel Cletheroe

Microsoft Research (United Kingdom) <https://orcid.org/0009-0003-6444-9149>

Doug Kelly

Microsoft Research

Christos Gkantsidis

Microsoft Research <https://orcid.org/0000-0002-6898-2368>

Michael Hansen

Microsoft Research

Paul Jeha

Microsoft Research

Kirill Kalinin

Microsoft Research

Jim Kleewein

Microsoft

Babak Rahmani

Microsoft Research

Saravan Rajmohan

Microsoft

Victor Rühle

Microsoft <https://orcid.org/0000-0002-8957-7628>

Jannes Gladrow

Microsoft Research

Francesca Parmigiani

Microsoft Research <https://orcid.org/0000-0001-7784-2829>

Hitesh Ballani

Microsoft Research <https://orcid.org/0000-0003-1573-3314>

Physical Sciences - Article

Keywords: generative AI, analog computing, diffusion models, flow matching, optical computing, neuromorphic computing, fixed-point compute paradigm

Posted Date: March 18th, 2026

DOI: <https://doi.org/10.21203/rs.3.rs-8919479/v1>

License:   This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

Additional Declarations: There is **NO** Competing Interest.

Analog Diffusion Models

Heiner Kremer^{*,1}, Fabian Falck^{*,1}, Jiaqi Chu^{*,1,†}, Grace Brennan¹, Burcu Canakci¹, James H. Clegg¹, Daniel Cletheroe¹, Douglas J. Kelly¹, Christos Gkantsidis¹, Michael Hansen², Paul Jeha¹, Kirill P. Kalinin¹, Jim Kleewein², Babak Rahmani¹, Saravan Rajmohani², Victor Rühle², Jannes Gladrow¹, Francesca Parmigiani^{1,‡}, & Hitesh Ballani¹

¹Microsoft Research, Cambridge, UK.

²Microsoft, Redmond, WA, USA.

*These authors contributed equally to this work.

†Correspondence to: Jiaqi.Chu@microsoft.com.

‡Correspondence to: Francesca.Parmigiani@microsoft.com.

Abstract

As generative artificial intelligence (GenAI) drives computational demands to unprecedented scales^{1,2}, digital hardware is approaching fundamental limits³. Analog and optical systems⁴⁻²⁰ promise orders-of-magnitude efficiency gains, but translating these to application-level gains is challenging due to the mismatch between hardware primitives and algorithmic requirements. Here, we introduce Analog Diffusion Models (ADMs) which implement diffusion inference with an implicit integration scheme, formulating each diffusion step as a fixed-point problem amenable for acceleration by efficient analog hardware. At the same time, training remains identical to that of conventional diffusion models, allowing adoption of established scalable training approaches with no additional overhead. We validate ADMs on analog hardware using three-dimensional optics with 2,304 programmable weights. On hardware, we generate two-dimensional distributions and latent-space distributions for MNIST, FashionMNIST, and ExtendedMNIST, demonstrating the feasibility of executing multi-layer diffusion processes entirely on noisy, non-traditional hardware. The current prototype reaches fixed-point convergence in 10–15 μs per diffusion step, with projections to nanosecond-scale convergence with miniaturization. In simulation, across multiple datasets, backbone architectures, and model sizes ranging from 32 million

to 13 billion parameters, ADMs match the sample quality of standard methods with up to $16\times$ fewer diffusion steps. Most importantly, they could achieve efficiency gains of more than $100\times$ at the application level without sacrificing generation quality, $100\times$ from hardware acceleration, and an additional $1-2\times$ from algorithmic improvement, highlighting the multiplicative benefit of hardware–algorithm co-design. Together, these results establish ADMs as a scalable and general, hardware-aligned framework for low-latency and energy-efficient generative modeling on analog computing platforms.

Keywords: generative AI, analog computing, diffusion models, flow matching, optical computing, neuromorphic computing, fixed-point compute paradigm

1 Generative artificial intelligence (GenAI) has emerged as one of the most trans-
2 formative technological developments of recent years, driving advances in language
3 modeling, visual media generation, and scientific discovery. Within this landscape, dif-
4 fusion models^{21,22} have become the dominant paradigm for generating continuous data
5 modalities, including images²³, audio²⁴, video^{25,26}, and scientific data^{27,28}. Despite
6 their empirical success, diffusion models typically require many sequential evalua-
7 tions to iteratively transform noise into data, making inference compute-, energy- and
8 latency-intensive.

9 Efforts to improve the efficiency of diffusion models have followed two largely inde-
10 pendent directions: alternative hardware designs and algorithmic advances. In this
11 context, particular emphasis is on the inference stage, which is estimated to account
12 for 90% of the energy consumed in commercial deployments²⁹.

13 Firstly, as digital electronics approach physical and economic limits, analog and
14 optical computing⁴⁻²⁰ have re-emerged as promising alternatives, offering orders-of-
15 magnitude improvements in latency and energy efficiency for compute intense matrix
16 multiplications central to AI applications. However, analog hardware is susceptible to
17 noise and often comes with a limited set of operations different from what is required
18 for today’s models which are optimized for digital accelerators. Therefore translating

19 these gains to the application level remains an open challenge. Existing optical imple-
20 mentations of diffusion models use diffractive optical layers and hybrid analog-digital
21 systems^{11,12}. This limits them to either linear architectures, due to the absence
22 of inter-layer nonlinearities¹¹, or to designs dominated by digital computations¹².
23 Previous analog electronic implementations using resistive and in-memory comput-
24 ing platforms have demonstrated accelerated diffusion inference^{13,14}. However, these
25 efforts rely on non-standard diffusion formulations and training workflows, preventing
26 the use of conventional models and established pipelines, and thereby limiting scal-
27 ability, generality and adoption. To our knowledge, no existing diffusion framework
28 simultaneously supports scalable multi-layer architectures, fully analog inference, and
29 compatibility with well-established training procedures, three properties that are key
30 to a general, efficient, and scalable diffusion system.

31 Secondly, substantial algorithmic efforts have focused on reducing the number of
32 sampling steps required for high-quality generation, either by distilling or rectifying
33 the model^{30,31} or improving sampling schemes³²⁻³⁴. In the latter case, most existing
34 approaches rely on low-order explicit discretizations of the underlying ordinary dif-
35 ferential equation (ODE), such as Euler, Heun, or midpoint schemes³³. While these
36 methods are computationally efficient, requiring only one or two function evalua-
37 tions per step, they can become unstable when the vector field varies rapidly with
38 its inputs. For diffusion this is the case particularly near the data manifold and close
39 to pure noise³⁵. In such regimes, explicit solvers often require increasingly small step
40 sizes to maintain stability and may produce overly smooth solutions when uncertainty
41 is high³³. Implicit solvers, which rely on implicit discretizations of the underlying
42 ODE, offer a promising alternative in these settings, but have seen limited adoption
43 in diffusion models since their iterative nature causes overheads when deployed on
44 conventional digital hardware. This motivates implementing implicit diffusion sam-
45 pling on analog hardware, where physical equilibration efficiently realizes the iterative

46 computations. The impact of analog noise, a central challenge of such non-traditional
47 accelerators, is mitigated both by the attractor dynamics of the fixed-point itera-
48 tions and by the inherent denoising nature of the diffusion process. Importantly, to
49 ensure that the resulting algorithms can be broadly used across models, datasets, and
50 hardware platforms, rather than being confined to specialized formulations, such an
51 approach must remain fully compatible with standard diffusion training pipelines. We
52 achieve this through careful co-design of the algorithms with the inherent properties
53 of the hardware, without sacrificing generality.

54 Here, we introduce Analog Diffusion Models (ADMs), a diffusion framework that
55 realizes diffusion steps through the physical equilibrium dynamics of analog hardware.
56 During inference, ADMs employ an implicit ODE solver, evaluating each diffusion
57 step via a fixed-point search. Crucially, training remains identical to that of con-
58 ventional diffusion models and does not require fixed-point search. This sets ADMs
59 apart from native fixed-point approaches such as deep equilibrium models³⁶ and their
60 diffusion variants^{37–39}, which require fixed-point search during training, increasing
61 computational cost and complexity and reducing generality.

62 The fixed-point search required by ADMs for inference can be implemented on
63 a range of analog platforms, either through energy-minimizing physical systems^{40,41}
64 or via architectures that implement iterative computations, involving matrix multi-
65 plications and nonlinear operations in a closed feedback loop^{20,42–44}. These iterative
66 architectures can be realized on analog hardware such as optical, resistive, memristive,
67 or neuromorphic systems, provided they support these core operations^{5,13,14,18,19,42,43}.
68 We adopt the iterative architecture²⁰ and perform analog diffusion inference on an
69 Analog Optical Computer (AOC). AOC performs fixed-point search entirely in the
70 analog domain by combining three-dimensional optics and analog electronics in a
71 loop until convergence (Fig. 1a). By avoiding repeated digital–analog conversions and

72 merging compute and memory, such non-traditional hardware can achieve substantial
73 efficiency gains of up to three orders of magnitude over digital accelerators^{5,18–20}.

74 With ADMs, training is performed on digital graphics processing units (GPU),
75 which conventional diffusion training is optimized for. To account for non-idealities
76 of the analog system, we train our models using a digital twin (DT) of the hard-
77 ware, which applies hardware specific alterations and corruptions to the parameters
78 and outputs of a model. The model itself is kept in standard form and the digital
79 twin can be seen merely as a differentiable translation layer to the simulated hardware
80 output for given standard neural network. Unlike some prior approaches that incor-
81 porate hardware-in-the-loop training or iterative digital-twin refinement within the
82 learning loop^{11,45–47}, we employ a single, task-agnostic DT whose non-model param-
83 eters are kept fixed across all experiments. This preserves standard diffusion training
84 while accounting for hardware non-idealities in the forward pass, and eliminates the
85 need for hardware evaluations during training. This also allows for a wide range of
86 hardware implementations, including those that are unsuitable for backpropagation
87 based training⁴⁵.

88 We benchmark ADMs across a range of experimental settings. First, we demon-
89 strate fully analog diffusion sampling on an AOC hardware with 2,304 optical weights
90 using diverse two-dimensional (2D) distributions (Fig. 1b, Fig. 2a). To our knowl-
91 edge, this constitutes the first demonstration of multi-layer diffusion models operating
92 entirely on scalable analog hardware while maintaining a standard training pipeline.
93 Second, we extend to larger image generation via a latent diffusion framework⁴⁸.
94 Executing the iterative diffusion process entirely on AOC and performing only the
95 final decoding digitally, we generate MNIST, FashionMNIST (FMNIST), and Extend-
96 edMNIST (EMNIST) samples (Fig. 1c). Third, we study the algorithmic benefit and
97 limitations of ADMs at scale using a large-scale DT and common backbone archi-
98 tectures, including convolutional U-Nets⁴⁹ and diffusion transformers (DiT)⁵⁰, on

99 high-resolution image datasets, such as Smithsonian Butterflies, AFHQ, and Oxford
100 Flowers-102. Across all architectures, datasets and model sizes, ranging from 32 mil-
101 lion to 13 billion parameters, ADMs and its variants consistently reduce the number
102 of diffusion steps required to reach a given sample quality. In particular, they achieve
103 the same Fréchet Inception Distance (FID) as conventional diffusion models with up
104 to $16\times$ fewer steps.

105 Importantly, ADMs translate the entire low-level, per-operation efficiency gains
106 of next-generation analog hardware into application-level improvements. By trading
107 off diffusion steps and fixed-point iterations, ADMs combine the demonstrated algo-
108 rithmic improvements with the projected analog hardware energy gains to indicate
109 a potential $>100\times$ improvement in energy efficiency, or equivalently latency, relative
110 to standard diffusion sampling running on digital hardware at the same generation
111 quality (Fig. 1d and Discussion). These results establish ADMs as a general, efficient,
112 scalable and low-latency route towards fully-analog generative modeling.

113 **Analog realization of fixed-point diffusion steps**

114 Diffusion models iteratively transform Gaussian noise into samples from a target dis-
115 tribution by integrating either a stochastic or an ordinary differential equation^{35,51}.
116 Here we focus on ODE-based diffusion variants, specifically the recently proposed
117 flow-matching formulation^{31,51,52}. Our method extends directly to other deterministic
118 sampling methods like Denoising Diffusion Implicit Models³².

119 While the iterative nature of ODE solvers is conceptually compatible with analog
120 hardware, direct implementation faces two mismatches. First, diffusion models inte-
121 grate their underlying ODE over a finite time horizon, typically from $t = 0$ to $t = 1$,
122 whereas analog systems evolve according to physical relaxation dynamics toward a
123 stationary state, effectively solving the ODE as $t \rightarrow \infty$. Second, early termination at
124 a precise time is generally difficult on analog systems: unlike digital implementations

125 with well-defined discrete iterations, many analog platforms^{13,14,20} evolve through
126 continuous-time dynamics without an intrinsic notion of iteration count. To address
127 these issues, we implement the diffusion process using an implicit ODE solver, where
128 each step corresponds to a fixed-point search. The choice of ODE solver affects the
129 number of steps required to reach a given sample quality^{33,53}. While conventional dif-
130 fusion solvers employ explicit discretizations of the underlying ODE (Eq. (3)), ADMs
131 rely on an implicit discretization scheme (Eq. (4)), requiring the solution of a nonlin-
132 ear system at each step. This formulation naturally aligns with equilibrating analog
133 hardware, defining each diffusion step as the stable fixed-point of a learned nonlinear
134 transformation conditioned on the previous state and diffusion time.

135 This formulation avoids the stability limitations of explicit discretizations in
136 rapidly varying vector field regimes near the data manifold and pure noise³⁵. It
137 further provides robustness to analog noise, arising from two complementary mecha-
138 nisms. First, the attractor dynamics of the fixed-point iteration make each diffusion
139 step inherently robust to noise (Supplementary Fig. 4a). Second, the diffusion pro-
140 cess itself further suppresses noise, since each denoising step is designed to operate on
141 noisy inputs, with additional analog noise effectively absorbed into the diffusion noise
142 (Supplementary Fig. 4b-c).

143 Model parameters are learned using standard diffusion objectives on digital
144 GPUs, without fixed-point search during training. Because analog non-idealities affect
145 inference-time behavior, in analog models training and inference are inherently cou-
146 pled: to reflect the non-idealities correctly, inference needs to be based on exactly the
147 same operations as training, which in this case would require prohibitively costly fixed-
148 point training. We avoid this by capturing these non-idealities in the DT and exploiting
149 a duality between explicit and implicit Euler discretization schemes. This ensures the
150 competing goal of consistency between training and inference while preserving stan-
151 dard training procedures, yet enables implicit ODE discretizations in ADMs that

152 are naturally suited to equilibrating analog hardware with inherent noise. Details on
153 methodology and implementation are provided in the Methods.

154 While ADMs can be implemented across diverse analog hardware systems, we
155 demonstrate their operation on AOC that integrates 3D optical and analog electronic
156 technologies in a recursive loop. We use micro light-emitting diode (microLED) arrays,
157 reflective spatial light modulators (SLMs) with pixels encoding model weights, and
158 photodetector (PD) arrays, leveraging fan-out and fan-in of light in the third dimen-
159 sion to enable inherently parallel, scalable matrix multiplications. Compared to the
160 system presented in ²⁰, our hardware is scaled up by a factor of about 10 \times , to 2,304
161 programmable optical weights, and has been re-designed (e.g. implementing two inde-
162 pendent optical paths instead of a single pass to represent positive and negative weights
163 simultaneously) to reduce noise, mitigate non-idealities, and minimize discrepancies
164 between experimental results and DT (refer to [Extended Data Fig. 1](#) and Methods
165 for additional hardware details). [Fig. 1a](#) describes the modular design for implement-
166 ing multi-layer diffusion models on non-traditional hardware: each layer of the neural
167 network is represented by a module, implementing optical matrix-multiplication and
168 analog non-linearity. Multiple modules are connected sequentially through a crossbar
169 analog switch, and the final module is connected back to the first to form a closed loop
170 for the implicit scheme. The same principle can be implemented using a single mod-
171 ule with one of the blockmatrix structures exemplarily shown in the inset. The latter
172 approach has been used in our experiments. Representative 2D distributions generated
173 on hardware are shown on the right. [Fig. 1b](#) provides a detailed visualization of the
174 diffusion process using hardware results. Each diffusion step consists of a fixed-point
175 search (columns). Once converged, the result of the diffusion step is used to condition
176 the next fixed-point search (rows). On this system, fixed-point inference for 2D distri-
177 butions takes 10–15 μs per diffusion step, governed by the current system bandwidth
178 of 5 MHz (-3dB; [Supplementary Fig. 1b](#)). Such per-step latencies are expected to be

179 on the order of nanoseconds in future hardware generations with target modulation
180 speeds in the tens of GHz^{4,54}.

181 Fig. 2a presents samples from a diverse set of 2D distributions generated on
182 the same hardware, showing that analog fixed-point search reliably implements dif-
183 fusion steps despite hardware noise and device nonidealities. Across the 2D target
184 distributions, the hardware results closely match the DT results, with a per-step nor-
185 malized root mean square error (NRMSE) around 10^{-2} (Supplementary Fig. 2c). For
186 quantitative benchmarking, we compare distributional discrepancies of the generated
187 distributions to the target distribution with those of a Gaussian distribution used as
188 a metric reference (Fig. 2b and c; *Supplementary Section D* for metric definitions).
189 In all cases, our results significantly improve over the Gaussian reference. To reduce
190 remaining stochastic errors, we average results either across multiple experimental
191 runs (Extended Data Fig. 7a) or over convergence windows (Extended Data Fig. 7b),
192 which yields stable outputs. Due to the noise robustness, five repetitions corresponding
193 to convergence window of $0.2 \mu\text{s}$ is sufficient.

194 For completeness, in addition to ADMs, our methodology enables the execution
195 of conventional (explicit) diffusion ODE steps on analog hardware (*Supplementary*
196 *Section C.2*). Explicit schemes are not governed by attractor dynamics unlike implicit
197 schemes, and are therefore more sensitive to analog noise. In the following we will
198 focus on the implicit version for its better hardware alignment and superior efficiency.

199 **Latent ADMs on AOC**

200 Following the state-of-the-art in diffusion modeling, we employ a latent diffusion
201 framework⁴⁸, performing the diffusion process in the latent space of a variational
202 autoencoder (VAE). At training time, an ADM learns the latent space distribution of
203 the VAE-encoded data. At inference time, the latent ADM is evaluated on the hard-
204 ware and the final result mapped to the data space using the VAE decoder (Fig. 3a).

205 The diffusion process runs entirely in the analog domain, while only the final decod-
206 ing step remains digital. In latent diffusion at scale, energy consumption is dominated
207 by the iterative diffusion process in latent space, whereas VAE decoding is a single
208 one-time operation, making the diffusion process the key bottleneck and particularly
209 well suited for analog acceleration. The all-analog nature of our diffusion process is
210 in contrast to prior optical diffusion work¹², where a large digital decoder network is
211 employed in every diffusion step, leaving only a small fraction of the diffusion process
212 to be implemented in the optical domain.

213 Using the latent ADM approach we generate samples from the canonical MNIST,
214 FMNIST and EMNIST datasets. Fig. 3b compares samples generated with ADM on
215 the AOC hardware with those obtained using ADM on the DT and a standalone VAE
216 (VAE-alone). Results show that ADMs (hardware) substantially outperform the VAE-
217 alone baselines qualitatively (Fig. 3b) and quantitatively, yielding lower FID in the
218 data space (Fig. 3c) and lower distributional discrepancies in the latent space (Fig. 3d).
219 Extended Data Fig. 8 provides additional representative results for EMNIST, showing
220 the evolution of the ADM (hardware) samples across diffusion steps and the robust-
221 ness to analog noise across trials and time intervals. Supplementary Fig. 5 further
222 shows close correspondence between the latent distributions generated by the hard-
223 ware, the DT and the ground truth on MNIST. These benchmarks demonstrate that
224 latent ADMs provide a hardware-compatible pathway for high-dimensional generative
225 modeling with our prototype hardware, enabling high-fidelity image synthesis while
226 confining all iterative diffusion computations to the analog domain.

227 **Analog Diffusion at Scale**

228 For benchmarking beyond the current hardware scale, we evaluate the implicit sam-
229 pling scheme of ADMs on multiple high-dimensional image datasets comprising
230 Smithsonian Butterflies, AFHQ, and Oxford Flowers-102, and across a range of model

231 architectures. These include convolutional U-Nets, DiT, and a scaled-up DT of our ana-
 232 log hardware without non-idealities. In addition to the implicit Euler scheme discussed
 233 so far (labeled as Implicit), we investigate the Crank-Nicolson solver (Implicit-CN),
 234 an implicit variant that combines the implicit update with an additional explicit
 235 model evaluation, enabling direct compatibility with our hardware (see Methods).
 236 We compare these approaches with the explicit Euler scheme (Explicit) employed in
 237 state-of-the-art models^{26,55}. We consider an image resolution of 256×256 and employ
 238 a standard latent diffusion approach using the frozen Stable Diffusion VAE²⁶ and
 239 models with sizes in the range between 32 million and 13 billion parameters. All our
 240 models are trained using the standard flow matching objective without loss weighting
 241 or rectification steps. For each experiment the same model is evaluated with differ-
 242 ent schedulers. Architectural details and training hyperparameters are provided in
 243 Supplementary Section A.3.

244 Across datasets and backbone architectures, implicit sampling consistently out-
 245 performs explicit sampling in the low-step regime, while converging to the same
 246 asymptotic sample quality as the number of diffusion steps increases (Fig. 4a). For
 247 example, on Oxford Flowers-102 using a U-Net backbone, the implicit sampler, and
 248 similarly Implicit-CN, achieves an FID of 26.5 using only 8 diffusion steps (with 4 fixed
 249 point iterations per step), outperforming the explicit variant with 128 steps (FID 30.6),
 250 corresponding to a $16\times$ reduction in steps and $4\times$ reduction in total neural network
 251 evaluations (algorithmic improvement), see middle graph in Fig. 4a. On Smithsonian
 252 Butterflies, using the DT of our hardware, both implicit variants achieve a smaller
 253 but consistent reduction in steps ($4\times$) compared to explicit, and up to $2\times$ algorithmic
 254 improvement, see left most graph in Fig. 4a. On AFHQ with DiT, the implicit sampler
 255 shows no advantage. This is explained by the stability of the fixed-point search. Stable
 256 convergence imposes a lower bound on the number of diffusion steps, T_{\min} , which scales
 257 with the Lipschitz constant of the learned vector field v , $T_{\min} \sim \text{Lip}(v)$, see Theorem

258 1 in Method. In the case of AFHQ with DiT, this constraint limits the effectiveness
 259 of basic implicit methods. Here, the Crank–Nicolson variant, Implicit-CN, provides a
 260 more favorable trade-off. As a second-order method with $\mathcal{O}(\Delta t^2)$ global error, com-
 261 pared to $\mathcal{O}(\Delta t)$ for explicit and implicit Euler, it improves accuracy at larger step
 262 sizes, while its fixed-point formulation admits a more relaxed stability condition, with
 263 $T_{\min, \text{CN}} \sim 2 \times \text{Lip}(v)$. On AFHQ, where the Implicit scheme is unstable at low diffusion
 264 step counts, Implicit-CN restores stable fixed-point convergence and recovers the effi-
 265 ciency gains of implicit sampling. This yields an up to $8\times$ reduction in diffusion steps,
 266 comparing 8 steps of Implicit-CN (FID 19.2) with 64 steps of Explicit (FID 18.7) and
 267 1-2 \times reduction in total neural network evaluations, see most right graph of Fig. 4a.

268 Fig. 4b shows qualitative samples comparing implicit variants (Implicit for Smith-
 269 sonian Butterflies and Oxford Flowers-102, and Implicit-CN for AFHQ) with the
 270 explicit Euler scheme at low diffusion step counts across datasets and architectures. For
 271 the same number of steps, implicit sampling yields sharper images with clearer struc-
 272 ture and greater sample diversity, whereas the explicit baseline exhibits noticeable blur
 273 and loss of detail. These qualitative differences are consistent with the corresponding
 274 improvements observed in FID.

275 Considering both, the number of diffusion steps and fixed-point iterations, com-
 276 bined with projected energy efficiencies of analog hardware, we investigate the
 277 potential energy efficiency of implicit solvers. Fig. 4c shows FID versus projected
 278 energy for different implicit solver configurations on the butterfly dataset compared
 279 to conventional explicit Euler sampling on digital hardware, each for 8, 16, 32 and 64
 280 diffusion steps respectively. A detailed examination of these projections is deferred to
 281 the Discussion.

282 Extended Data Fig. 2 shows that implicit samplers approach convergence with as
 283 little as two fixed-point iterations; more iterations can improve the results especially
 284 for small numbers of diffusion steps but show diminishing improvement for 8 and

285 more steps. In particular for eight or more diffusion steps implicit sampling with two
286 iterations yields similar or better results than a higher order Heun sampler, which is
287 a commonly used best practice³³ and uses the same number of function evaluations.
288 Additional qualitative results are presented in [Extended Data Fig. 3](#), [Extended Data](#)
289 [Fig. 5](#), [Extended Data Fig. 4](#).

290 Finally, [Fig. 4d](#) and [Extended Data Fig. 6](#) qualitatively show that the implicit
291 paradigm can be beneficial even for highly optimized state-of-the art models. We use
292 FLUX.1⁵⁵, a 13 billion parameter text-to-image diffusion model. As a rectified flow
293 model optimized for few-step explicit Euler sampling, the model produces high quality
294 results already at 5 explicit steps. To explore a possible advantage of implicit meth-
295 ods we thus reduce the number of steps to 2. The standard implicit scheme does not
296 converge at such low step counts, as per previous discussion. However, Implicit-CN
297 achieves robust fixed-point convergence: with 2 diffusion steps and 4 fixed-point itera-
298 tions, it significantly improves sample quality, producing sharper fine-grained details.
299 These initial findings underline the generality of the approach, extending beyond
300 controlled experiments to industry-grade diffusion models.

301 Taken together, our results demonstrate that implicit diffusion samplers constitute
302 a general and practically effective alternative to explicit schemes, with clear benefits
303 that persist from medium-scale image models to highly optimized frontier systems.

304 Discussion

305 In this work, we introduce ADMs, a diffusion framework that co-designs implicit
306 solvers with non-traditional hardware to address three key challenges: generality,
307 efficiency, and scalability, for diffusion atop non-traditional hardware.

308 By confining fixed-point search to inference while retaining standard diffusion
309 training, ADMs overcome the training challenges typically associated with fixed-
310 point models such as deep equilibrium models, while harnessing the energy efficiency

311 and noise robustness of analog fixed-point dynamics, which complement the noise-
312 suppressing properties of diffusion denoising. This ensures broad applicability and
313 generality of ADMs.

314 We demonstrate the feasibility of ADMs on a 2,304-optical-weight AOC, gener-
315 ating diverse 2D distributions, and extending the approach via latent diffusion to
316 enable implicit sampling on AOC for benchmarking across several larger datasets. In
317 digital simulations, we benchmark ADMs across a range of high-dimensional image
318 generation tasks, model backbones and sizes. Across all settings, implicit variants con-
319 sistently reduce diffusion steps by up to $16\times$ and network evaluations by $1-2\times$ while
320 maintaining asymptotic sampling quality. While implicit sampling allows using fewer
321 diffusion steps, its reliance on fixed-point convergence imposes a minimum number
322 of diffusion steps dependent on the backbone model’s Lipschitz constant. This can
323 lead to failure modes for the implicit Euler method, as observed on AFHQ using
324 DiT and FLUX. The Crank-Nicolson method (Implicit-CN), a hardware compatible
325 implicit variant that combines higher-order accuracy with relaxed convergence con-
326 straints, overcomes this limitation, consistently producing high-quality samples with
327 stable convergence, and achieves an algorithmic reduction of required function eval-
328 uations of approximately $2\times$. These results position ADMs as a promising pathway
329 toward more compute-efficient diffusion sampling, even atop digital hardware, and
330 motivate further investigation.

331 The main advantage of ADMs is that they retain operation-level efficiencies of
332 analog hardware at the application level and further compound these with algorithmic
333 improvements. Estimated analog efficiencies reach up to femtojoule-per-operation or
334 peta operations per second per watt^{5,18–20}, suggesting potential energy savings of
335 $100\times$ to $1000\times$ relative to current digital accelerators. Assuming a conservative $100\times$
336 per-operation energy improvement and the $1-2\times$ algorithmic advantage achieved with
337 our proposed hardware model architecture, ADMs yield projected energy savings,

338 or equivalent latency reductions, of more than $100\times$ versus explicit digital sampling
339 (Fig. 4d). These gains arise from the co-design of our method with hardware properties,
340 effectively translating low-level analog efficiency into application-level improvements.

341 Realizing these gains for industrial and scientific applications requires analog hard-
342 ware that matches the scale and precision requirements of modern generative models.
343 Advances in low-bit and quantized inference on digital hardware⁵⁶ are expected to
344 carry over to analog implementations, further increasing their efficiency. Scaling is
345 particularly critical for high-dimensional image generation tasks, which demand mil-
346 lions or billions of weights. Significant progress has recently been demonstrated in this
347 direction^{4,7}, but planar optical approaches^{4,7,10,18} remain limited by chip reticle size,
348 as routing and computation share the same plane. In contrast, the AOC hardware we
349 used for the ADM implementation leverages the third dimension to separate compu-
350 tation from data routing, enabling improved scalability. Each SLM pixel can serve as
351 an individual weight modulator, with a footprint that is several orders of magnitude
352 smaller than in planar alternatives. Off-the-shelf SLMs support up to four million μm -
353 size pixel arrays allowing to represent matrices of sizes up to 1000×1000 per module,
354 while emerging metasurface arrays promise 9.5 million pixel with microsecond-scale
355 refresh rates⁵⁴, opening new scaling regimes for optical accelerators. Although cur-
356 rent 3D prototypes rely on chip-sized technologies and bulky systems^{12,19,20}, future
357 platforms will require full integration of the optical stack within centimeter-scale vol-
358 umes. The modular design in Fig. 1a enables layers of up to four million parameters
359 with current SLMs, comparable to the feedforward layers in modern models⁵⁵. Larger
360 layers can be achieved by distributing the optical matrix-vector multiplications over
361 replicated modules. Finally, realizing full hardware efficiency requires analog memory,
362 e.g., a sample-and-hold circuit, to avoid digital conversions between diffusion steps.

363 Overall, the results presented here point toward a new trajectory for GenAI. By
364 grounding diffusion inference in the relaxation of a physical system rather than digital

365 iteration, ADMs offer a general and scalable route to high-throughput, energy-efficient
366 and low-latency generative modeling. As analog accelerators mature and expand in
367 scale, this hardware–algorithm alignment will be a promising path to reduce the
368 computational and environmental burden of GenAI workloads.

369 Methods

370 Flow matching and definitions

371 Sampling from a flow-matching based diffusion model corresponds to integrating the
372 ODE

$$\frac{d}{dt}x(t) = f_\theta(x(t), t), \quad x(0) = x_0, \quad (1)$$

373 from $t = 0$ to $t = 1$, where $x_0 \sim N(0, I)$ is a sample from a multi-variate Gaussian
374 and in our case $f_\theta : \mathbb{R}^d \times \mathbb{R}_+ \rightarrow \mathbb{R}^d$ denotes a time-conditional, K -layer feedforward
375 neural network of the form

$$f_\theta(x, t) = (L_K \circ \dots \circ L_1)(x, t), \quad L_i(x, t) = W_i \sigma(x) + b_i + t e_i. \quad (2)$$

376 Each layer is a map $L_i : \mathbb{R}^{d_{i-1}} \times \mathbb{R}_+ \rightarrow \mathbb{R}^{d_i}$, with weight $W_i \in \mathbb{R}^{d_i \times d_{i-1}}$, bias $b_i \in \mathbb{R}^{d_i}$,
377 and time embedding vector $e_i \in \mathbb{R}^{d_i}$, $i = 1, \dots, K$, and $d_0 = d_K = d$. The trainable
378 model parameters are $\theta := \{W_1, b_1, e_1, \dots, W_K, b_K, e_K\}$. A common ODE solver is
379 the explicit Euler discretization, which for a fixed number of steps $T = \frac{1}{\Delta t}$ iterates
380 the update

$$x_{t+\Delta t} = x_t + \Delta t f_\theta(x_t, t). \quad (3)$$

381 In contrast the implicit Euler scheme uses the implicit update rule

$$x_{t+\Delta t} = x_t + \Delta t f_\theta(x_{t+\Delta t}, t + \Delta t), \quad (4)$$

382 which, as $x_{t+\Delta t}$ appears on both sides, requires solving a system of non-linear equations
383 at each step. Importantly, both updates use the same operator $\text{Euler}_{x_t, t, \Delta t, \theta}(z) :=$
384 $x_t + \Delta t f(z, t; \theta)$ merely evaluated at different $t \in [0, 1]$ and $z \in \mathbb{R}^d$.

385 **ADMs: Methodology**

386 The idea behind ADMs is to construct a hardware compatible operator, which can
387 be evaluated both explicitly during training (Eq. (3)) and implicitly during inference
388 (Eq. (4)). In the following we will derive this construction.

389 ADMs are implemented by embedding the network (2) into a higher-dimensional
390 analog state space \mathbb{R}^n and realizing its computation through repeated application of
391 a fixed update rule. We denote by

$$\text{ADM}_{x,t;\theta}(z) = W\sigma(z) + b + te + \text{pad}(x) \quad (5)$$

392 the ADM operator, a re-parameterization of the update implemented by AOC, where
393 $z \in \mathbb{R}^n$ is the internal AOC state, $b, e \in \mathbb{R}^n$ are concatenations of biases b_i and
394 time embeddings e_i respectively, and $\text{pad}(x) := (x, 0, \dots, 0)^\top \in \mathbb{R}^n$ embeds the dif-
395 fusion state $x \in \mathbb{R}^d$ into the analog state space \mathbb{R}^n . AOC iterates this update until
396 convergence to a fixed-point $z^* \in \mathbb{R}^n$, with $\text{ADM}_{x,t;\theta}(z^*) = z^*$.

397 The purpose of the digital twin is to incorporate the non-idealities expected at
398 inference time during training. However, this makes it crucial that training and infer-
399 ence rely on the same update operator. The below construction of the ADM operator
400 ensures this property.

401 ***Blockmatrix layout***

402 Multi-layer ADMs can be implemented either combining multiple AOC systems in a
403 modular way (Fig. 1a), where all modules participate in a shared recurrent loop, or by
404 using a single AOC system with the block structure shown in the bottom of Fig. 1a
405 (implicit). Both implementations are mathematically equivalent and we will focus on
406 the latter in the following.

407 The block-structured matrix W and bias b of ADMs are constructed such that
 408 successive applications of the ADM operator correspond to cycling the signal through
 409 the layers of the network. Specifically, with each application of the ADM operator the
 410 signal gets propagated down the state and the output of the last layer gets fed back
 411 to the top. For $\Delta t = \frac{1}{T}$ we define $\theta(\Delta t) = \{W_{\Delta t}, b_{\Delta t}, e_{\Delta t}\}$ with

$$W_{\Delta t} = \begin{bmatrix} 0 & 0 & \cdots & 0 & \Delta t W_K \\ W_1 & 0 & \cdots & 0 & 0 \\ 0 & W_2 & \ddots & 0 & 0 \\ \vdots & \ddots & \ddots & 0 & \vdots \\ 0 & 0 & \cdots & W_{K-1} & 0 \end{bmatrix}, \quad b_{\Delta t} = \begin{bmatrix} \Delta t b_K \\ b_1 \\ b_2 \\ \vdots \\ b_{K-1} \end{bmatrix}, \quad e_{\Delta t} = \begin{bmatrix} \Delta t e_K \\ e_1 \\ e_2 \\ \vdots \\ e_{K-1} \end{bmatrix}, \quad (6)$$

412 as well as $\text{pad}(x) = (x, 0, \dots, 0)^\top \in \mathbb{R}^n$.

413 **Training**

414 Setting $\Delta t = 1$ at training time, after K applications of the ADM operator, the first d
 415 components (defined as $[\cdot]_{:d}$) of the AOC state recover the network output. In contrast
 416 to evaluation in the hardware which always iterates the operator until convergence to
 417 a fixed-point, during training on GPU using the DT we can evaluate the system after
 418 a finite number of ADM iterations. This means at training time we can evaluate the
 419 model explicitly via

$$f_\theta(x, t) = \left[(\text{ADM}_{x=0, t; \theta(1)})^K(x) \right]_{:d}, \quad (7)$$

420 where ADM^K denotes K applications of the ADM operator, and thus use standard
 421 (explicit) training procedures.

422 Using this formulation, model parameters θ are learned using standard conditional
 423 flow matching with optimal transport path⁵¹. Given samples $x_0 \sim p_0$ (Gaussian base
 424 distribution) and $x_1 \sim p_1$ (data distribution), as well as timestep $t \sim \text{Uniform}([0, 1])$

425 and the linear interpolation $x_t = tx_1 + (1 - t)x_0$, the training objective is

$$\mathcal{L}_{\text{CFM}} = \mathbb{E}_{(x_0, x_1), t} \|f_\theta(x_t, t) - (x_1 - x_0)\|_2^2, \quad (8)$$

426 with f_θ as defined in (7).

427 *Inference*

428 To correctly account for the non-idealities of the hardware, at inference time we can
 429 only apply the same operator learned during training. By construction of the ADM
 430 operator, K applications of $\text{ADM}_{x_t, t; \theta(\Delta t)}$ on x_t yield the explicit Euler update (3),
 431 i.e.,

$$\left[\text{ADM}_{x_t, t; \theta(\Delta t)}^K(x_t) \right]_{:d} = \text{Euler}_{x_t, t, \Delta t; \theta}(x_t) = x_t + \Delta t f_\theta(x_t, t). \quad (9)$$

432 The underlying working principle of ADMs is that iterating this update rule (with
 433 $t \rightarrow t + \Delta t$) beyond K iterations is still meaningful as the fixed-point of the Euler
 434 operator provides the implicit Euler update (4) and as the Euler operator itself can be
 435 represented by K ADM iterations, consequently the fixed-point of the ADM operator
 436 can be identified with the implicit Euler update. This is formalized by the following
 437 theorem.

438 **Theorem 1** (ADM Inference) *Let σ be $\text{Lip}(\sigma)$ -Lipschitz and let $\|\cdot\|_\infty$ denote*
 439 *the induced matrix norm. If $\Delta t < \frac{1}{\text{Lip}(\sigma)^K \|W\|_\infty^K}$, then for any (x, t) the operator*
 440 *$\text{ADM}_{x, t; \theta(\Delta t)}$ has a unique fixed point $z^* \in \mathbb{R}^n$ and, for any initialization z_0 , we have*
 441 *$\lim_{k \rightarrow \infty} \text{ADM}_{x, t; \theta(\Delta t)}^k(z_0) = z^*$. Moreover, the readout components $[z^*]_{:d}$ of the fixed-point*
 442 *solve the implicit Euler equation (4).*

443 Details and proofs can be found in *Supplementary Section C*.

444 Finally, combining these aspects yields our training and inference process described
 445 in Algorithm 1.

446 The same duality between explicit and implicit Euler scheme can be directly
 447 applied to make the Crank-Nicolson solver, whose update rule is given by $x_{t+\Delta t} =$
 448 $x_t + \Delta t/2 (f_\theta(x_t, t) + f_\theta(x_{t+\Delta t}, t + \Delta t))$, compatible with analog hardware. To this
 449 end, first evaluate $\tilde{x}_t = x_t + \Delta t/2 f_\theta(x_t, t)$ either in digital or using our analog
 450 version of explicit diffusion sampling (see *Supplementary Section C.2*). Using this pre-
 451 computed \tilde{x}_t , the result of the Crank-Nicolson scheme is given by the fixed-point of the
 452 operator $\text{ADM}_{\tilde{x}_t, t+\Delta t; \theta(\Delta t/2)}$, which is immediately compatible with AOC. Note that
 453 convergence of the Crank-Nicolson solver is already guaranteed for the less restrictive
 454 condition $\Delta t < \frac{2}{\text{Lip}(\sigma)^K \|W\|_\infty^K}$, which translates to a $2\times$ smaller minimum number of
 455 steps and allows more robust few-step sampling as demonstrated by the results on
 456 AFHQ in Fig. 4.

457 Relaxed fixed-point iterations to address hardware constraints

458 Optical vector-matrix multiplication (OVMM) requires fanning-out and fanning-in
 459 optical signals to and from all the modulators, i.e., the pixels of the SLM in AOC,
 460 representing the weights. As matrices become increasingly sparse, the total reflected
 461 or transmitted optical power decreases, making the system progressively more lossy.
 462 In the extreme case of an identity matrix, all optical power is attenuated except at
 463 a single modulator per row and column, rendering such operations impractical in
 464 the optical domain. To solve this general optical hardware limitation, we propose a
 465 novel methodology that is implemented in AOC as follows. We implement an analog
 466 skip connection that bypasses the OVMM for a fraction α of the signal, while the
 467 remaining $(1-\alpha)$ fraction is processed through the lossy optical path. During training,
 468 such a skip connection would break the duality between explicit and implicit Euler
 469 schemes that underlies ADMs. However, when applied as a post-processing step to
 470 a trained model, this mechanism corresponds to a relaxed fixed-point iteration and
 471 can be exploited without modifying the training procedure. Specifically, we define an

Algorithm 1 Training and inference of ADMs

Require: Network parameters $\theta_{\Delta t}$, step size Δt , total diffusion steps T

Require: (Digital twin of) ADM operator $\text{ADM}_{x,t;\theta}$

1: **Training (explicit evaluation using digital twin)**

2: **for** each minibatch **do**

3: Sample $x_0 \sim p_0$ (Gaussian), $x_1 \sim p_1$ (data)

4: Sample $t \sim \mathcal{U}[0, 1]$

5: Interpolate $x_t \leftarrow tx_1 + (1-t)x_0$

6: Evaluate network explicitly:

▷ K -layer model

$$f_\theta(x_t, t) \leftarrow \left[(\text{ADM}_{x=0,t;\theta_{\Delta t=1}})^K(x_t) \right]_{:d}$$

▷ $[\cdot]_{-d}$: for explicit

7: Compute loss

▷ Standard CFM loss

$$\mathcal{L} \leftarrow \|f_\theta(x_t, t) - (x_1 - x_0)\|_2^2$$

8: Update parameters θ via backpropagation and gradient descent

9: **end for**

10: **Inference (implicit evaluation using hardware):**

11: Initialize $x_0 \sim p_0$

12: **for** $n = 1, \dots, T$ **do**

13: Set $t_n \leftarrow n/T$

▷ $t_n \leftarrow (n-1)/T$ for explicit

14: Initialize AOC state $z^{(0)}$ arbitrarily

15: **repeat**

16: $z^{(k+1)} \leftarrow \text{ADM}_{x_n, t_n; \theta_{\Delta t}}(z^{(k)})$

▷ Fixed-point search

17: **until** convergence to fixed point z^*

▷ Theorem 1 guarantees convergence

18: Extract denoised update:

▷ Take first d entries ($[\cdot]_{-d}$: for explicit)

$$x_{n+1} \leftarrow [z^*]_{:d}$$

19: **end for**

20: **return** x_T

472 α -relaxed fixed-point iteration as

$$\text{ADM}_{x,t;\theta}^{(\alpha)}(z) := \alpha \odot z + (1 - \alpha) \odot \text{ADM}_{x,t;\theta}(z), \quad 0 \leq \alpha < 1, \quad (10)$$

473 which introduces an analog skip connection while preserving the structure of the ADM

474 operator. Crucially, after convergence the relaxed operator provides the same result as

475 the original ADM operator, since the mappings $x \mapsto f(x)$ and $x \mapsto \alpha x + (1 - \alpha)f(x)$

476 admit identical fixed-points. Under the same step-size conditions as in the non-relaxed

477 case, the α -relaxed operator converges to a unique fixed point corresponding to the

478 diffusion update (Theorem 6). The skip connection can therefore be introduced post-
479 training through a rescaling of the learned model parameters and only increases the
480 number of fixed-point iterations required for convergence.

481 By selecting α appropriately, either globally or in a channel-wise (per analog state
482 variable) manner, we can trade off optical power loss against the number of fixed-point
483 iterations. This enables the representation of large-weight matrices within the dynamic
484 range of the hardware while maintaining stable convergence and an improved effective
485 signal-to-noise ratio. Further details are provided in *Supplementary Section C.3*.

486 **Analog hardware platform**

487 Experiments were conducted on an analog optical computing platform comprising
488 2,304 programmable weights implemented through 3D optical matrix multipliers inte-
489 grated with channel-wise analog electronic nonlinearities in a recursive architecture
490 (see [Extended Data Fig. 1](#) for further details). During each fixed-point iteration, the
491 system state vector z_k switches back and forth between optical and analog electronic
492 forms. Matrix multiplications occur in the optical domain, where z_k is encoded as
493 light intensities from microLED arrays (peak wavelengths 503 nm at 20 mA), and
494 the weight matrix W is encoded onto reflective SLM pixels. Light emitted from each
495 microLED fans out across a row of the SLM for element-wise multiplication, with
496 resulting signals summed column-wise by a PD array as electrical voltages⁵⁷. Building
497 upon the optical matrix multiplication scheme described in²⁰, we implement separate
498 optical paths for positive and negative weights, each using an independent SLM, to
499 reduce system crosstalk. This scheme leverages 3D optics, with efficient fan-in and
500 fan-out enabled by spherical and cylindrical lenses, enabling inherently parallel and
501 scalable large-matrix multiplications. Analog voltages are processed through channel-
502 wise nonlinear activation functions (hyperbolic tangent), implemented with bipolar
503 differential pairs. The custom electronic board design is detailed in²⁰ and Extended
504 Data Figure 3 therein. Each analog iteration involve three steps: nonlinear activation,

505 matrix multiplication, and summation with fixed bias, time-embedding, and condition-
506 ing voltages. The same weights are reused at every iteration, following an in-memory
507 computing style. The system runs until the signals stabilize, reaching a fixed point
508 for each respective diffusion step. Finally, the voltages are recorded using a 48-port
509 digitizer (Spectrum M2p.5913).

510 **Sampling and experiment workflow**

511 Waveforms were measured with a sampling rate of 4.6875 MHz. Each experiment
512 captured 90 samples, corresponding to a total acquisition window of 19.2 μs . The
513 fixed-point search remained entirely in the analog domain during this interval. Sam-
514 ples were collected between 68% and 93% of the solve time, which translates to a time
515 window from 13.1 μs to 19.7 μs . This range ensures that equilibrium dynamics are
516 captured after initial transients but before full convergence. Once sampled, readout
517 voltages were linearly converted to solution variables. Voltages were averaged over 100
518 experiment for figures in the main text; averaging over 5 experiments shows compara-
519 tive generation quality (see [Extended Data Fig. 7](#) for 2D distributions and [Extended](#)
520 [Data Fig. 8c-d](#) for latent tasks).

521 **DT calibration procedure**

522 The DT is calibrated in two stages: an open-loop calibration that identifies component-
523 level hardware parameters, and a closed-loop calibration that refines system-level
524 dynamics. In the open-loop stage, we measure the behavior of each physical
525 block—electrical offsets and noise characteristics, the analog nonlinearity, microLED
526 input efficiency, SLM response, photodiode gain and crosstalk, and polarization-beam-
527 splitter leakage—using dedicated experiments designed to isolate each contribution.
528 Together, these component-level effects comprise the principal hardware nonidealities
529 discussed in the ADM methodology; characterizing them enables defining the initial
530 parametric form of the DT. To capture remaining linear distortions, we perform an

531 open-loop OVMM sweep over 100 random weight matrices, 6 input vectors, and 7
532 values of the gradient gain β (which scales the optical signal applied to the $W\sigma(z)$
533 term in (5)), and fit output-correction vectors. Across 4200 open-loop experiments, the
534 mean NRMSE was 0.015 (Supplementary Fig. 2a). This modular calibration workflow
535 is general in structure and is intended to be adaptable, in principle, to other analog
536 computing platforms comprising linear-nonlinear building blocks.

537 In the closed-loop stage, we validate and adjust the DT under recurrent operation.
538 Using 1000 experiments that span random weight matrices, channel-wise nonlinear
539 coefficients, and realistic β configurations, we refine a small set of global scaling param-
540 eters: the gains on the gradient, input, and the α skip connection term described
541 in (10), along with residual offsets. This yields a digital model that accurately
542 reproduces the hardware’s equilibrium trajectories across configurations used in the
543 diffusion experiments. Mean NRMSE across 1000 closed-loop experiments was 0.009
544 (Supplementary Fig. 2b).

545 Acknowledgments

546 We acknowledge colleagues at Microsoft Research, Cambridge, UK for fruitful
547 discussions.

548 Author contributions

549 H.K., F.F., and J.C. contributed equally to this work. J.C. and F.P. are corresponding
550 authors. J.G., F.P. and H.B. conceived and managed the project. H.K. and F.F. devel-
551 oped the framework and performed the hardware simulations for the 2D experiments.
552 H.K. derived the algorithm and theoretical results and performed the large scale sim-
553 ulations. F.F. performed the hardware simulations for the latent experiments. J.C.,
554 J.H.C., D.J.K., D.C., G.B. and F.P. proposed, designed and implemented the optical

555 and electrical parts. J.C. performed the experiments, analysed their data and imple-
556 mented the digital twin with support from J.H.C., D.J.K., D.C., G.B. and F.P. D.J.K.,
557 J.H.C., F.P., K.P.K. and D.C. developed the scaling roadmap. J.C., D.J.K., J.H.C.,
558 developed the control software. B.R., K.P.K., F.P., C.G., P.J., B.C., M.H., S.R., V.R.,
559 J.K. and H.B. provided technical input to the simulations and experiments. H.K.,
560 J.C., F.P., F.F., J.H.C. and K.P.K wrote the paper with input from all authors.

561 **Competing interests**

562 The authors of the paper have filed several patents relating to the subject matter
563 contained in this paper in the name of Microsoft Co.

564 **Code availability**

565 The code used in this study is currently being prepared for public release and will
566 be made available upon publication. During the review process, the code is available
567 from the corresponding author upon reasonable request.

568 **Data availability**

569 The data that support the findings of this study are currently being prepared for
570 public release and will be made available upon publication. During the review process,
571 the data are available from the corresponding author upon reasonable request.

572 **Additional information**

573 **Correspondence and requests for materials** should be addressed to Jiaqi Chu
574 or Francesca Parmigiani.

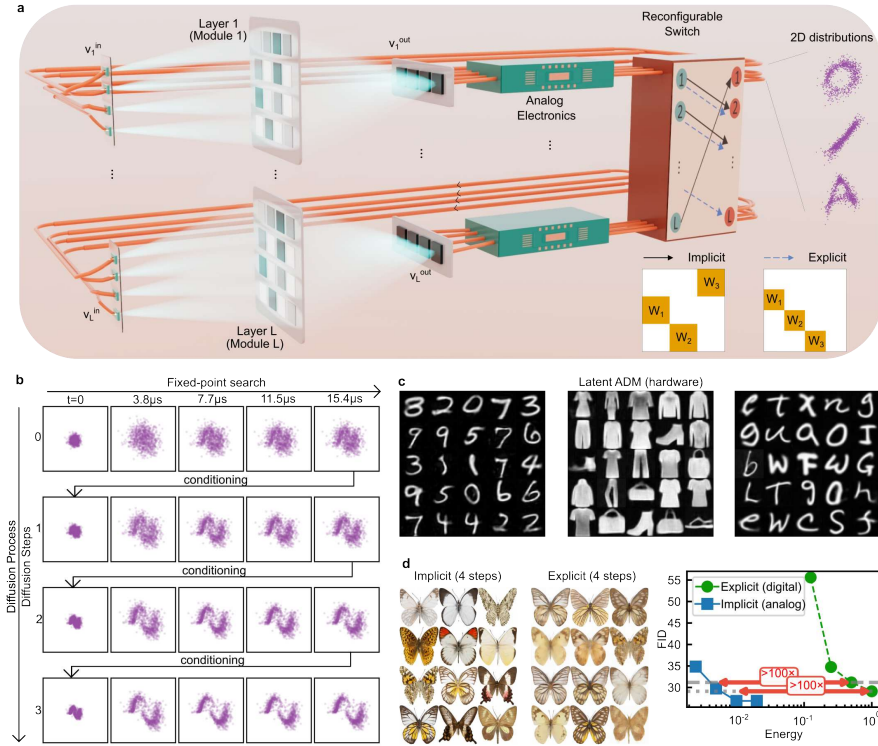


Fig. 1: Analog Diffusion Models (ADMs). **a**, ADMs implemented on modular analog optical computer (AOC) units. Multiple units are connected to implement the multi-layer structure of the diffusion model. A programmable switch selects between implicit and explicit evaluation modes. Mathematically, the multi-unit architecture is equivalent to a single system with a block-matrix layout, as illustrated in the inset. Each unit comprises optical vector–matrix multipliers implemented using a three-dimensional optical system with spatial light modulators, with nonlinearities provided by analog electronics. Conditioning information is injected into the loop as fixed voltages, and the recursive system converges to a fixed-point. Representative two-dimensional distributions generated by the hardware are shown on the right. **b**, The ADM sampling process proceeds along two axes: each diffusion step (rows) is implemented as a fixed-point search (columns) conditioned on the previous result (results on hardware). **c**, Selected latent ADM samples generated on the hardware for MNIST, FMNIST and EMNIST. **d**, Simulation of Implicit (analog) versus Explicit (digital) on the Smithsonian Butterflies dataset using the digital twin (DT) of the hardware. Qualitative images at 4 diffusion steps and quantitative FID versus projected energy for the two solvers, illustrating the excess of $100\times$ application-level energy efficiency gain.

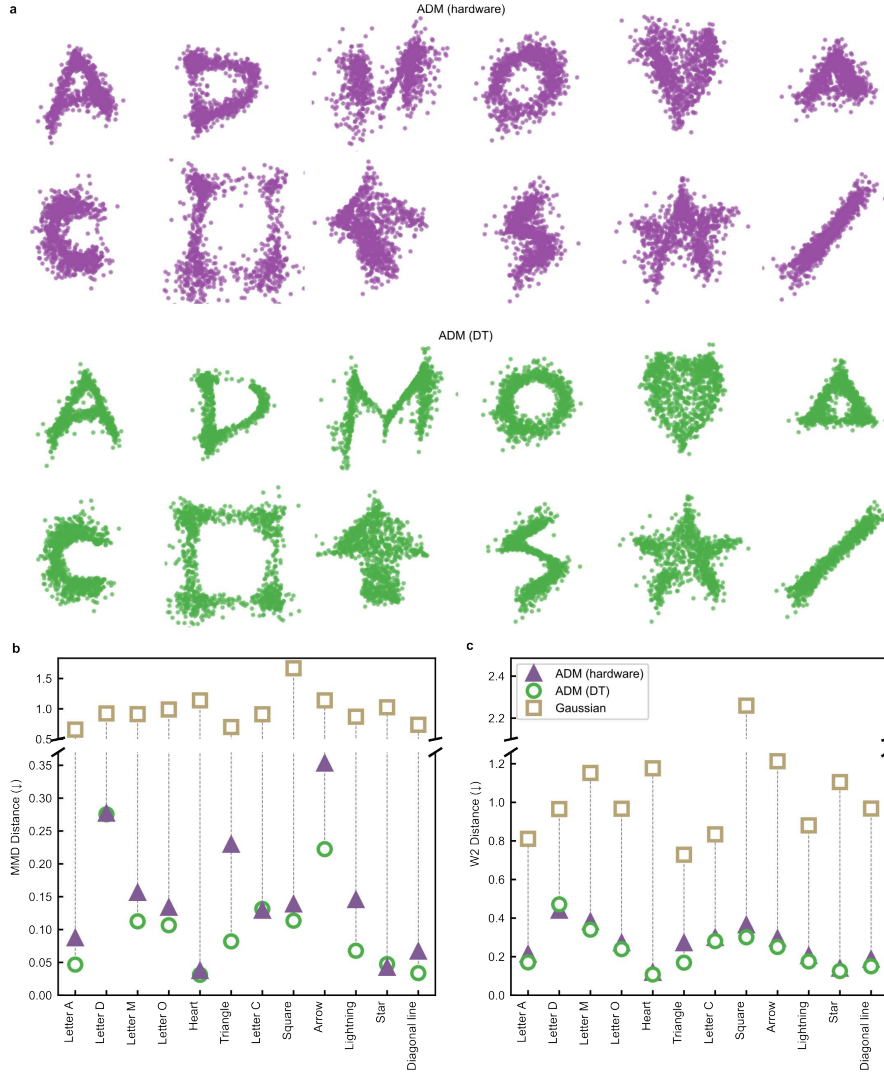


Fig. 2: ADM inference for 2D distributions. **a**, Examples of 2D distributions generated by ADMs on the hardware (top two rows) and its DT (bottom two rows). From left to right, we show: A, D, M, O, heart, triangle (top row), C, square, arrow, lightning, star and diagonal line (bottom row). **b**, Maximum Mean Discrepancy (MMD) between the ground truth and distributions generated by the hardware, the DT and a Gaussian reference distribution. The reference distribution is obtained by fitting a 2D multivariate Gaussian to the dataset points. All 2D point clouds are normalized to the square $[-1, 1]^2$ using the dataset distribution for each shape as reference. **c**, Wasserstein-2 (W2) distance between the groundtruth and generated distributions for multiple 2D datasets.

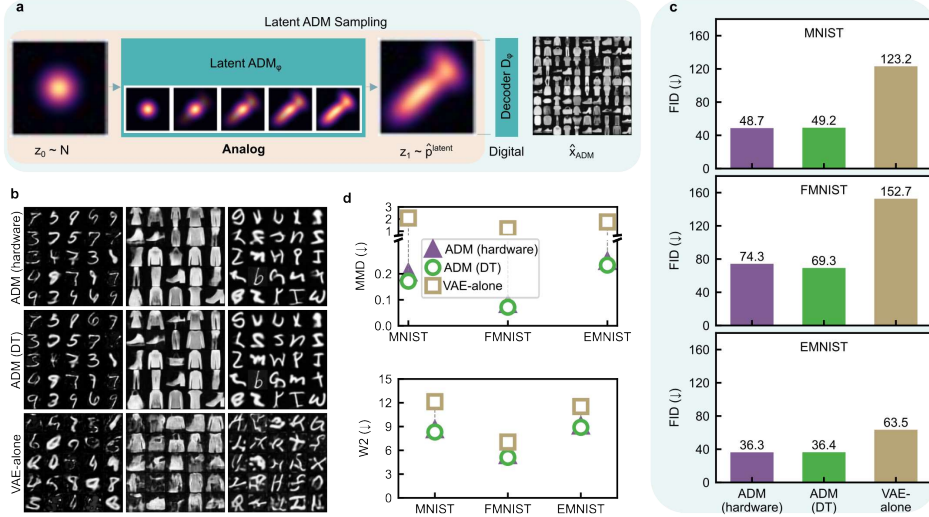


Fig. 3: Latent ADMs: operating principle, characterization and results. **a**, Schematic of the latent ADM framework. All diffusion steps are fully computed in analog hardware, and only the final decoding step is performed digitally. **b**, Samples generated by the latent ADM for MNIST (left), FMNIST (middle) and EMNIST (right), comparing ADMs on hardware (top), ADMs on the DT (middle) and the VAE-alone samples (bottom). **c**, Fréchet Inception Distance (FID) for MNIST (top), FMNIST (middle) and EMNIST (bottom), computed using 1,000 generated samples, comparing ADMs on hardware, ADMs on the DT, and the VAE-alone, highlighting the contribution of the ADM to the image generation. **d**, Latent-space distribution similarity between generated and target data (MMD and W2). Hardware results (triangles) closely match DT (circles), while VAE-alone (squares) serves as a baseline.

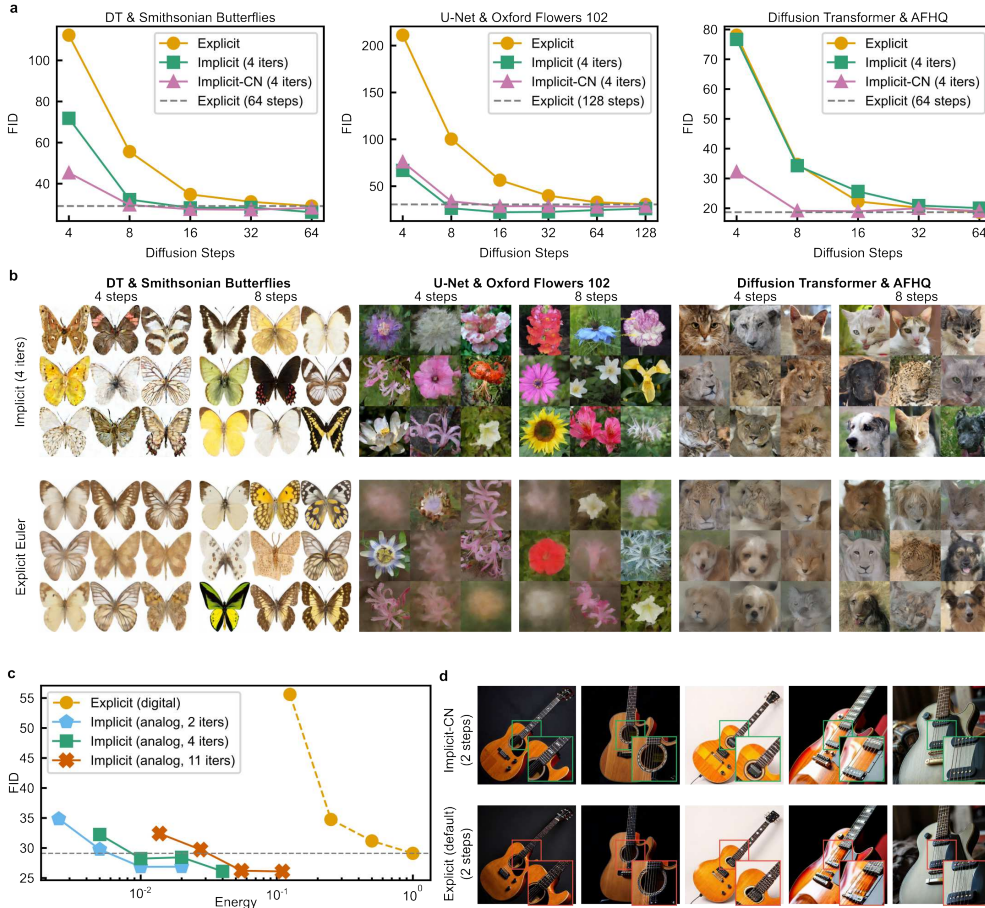
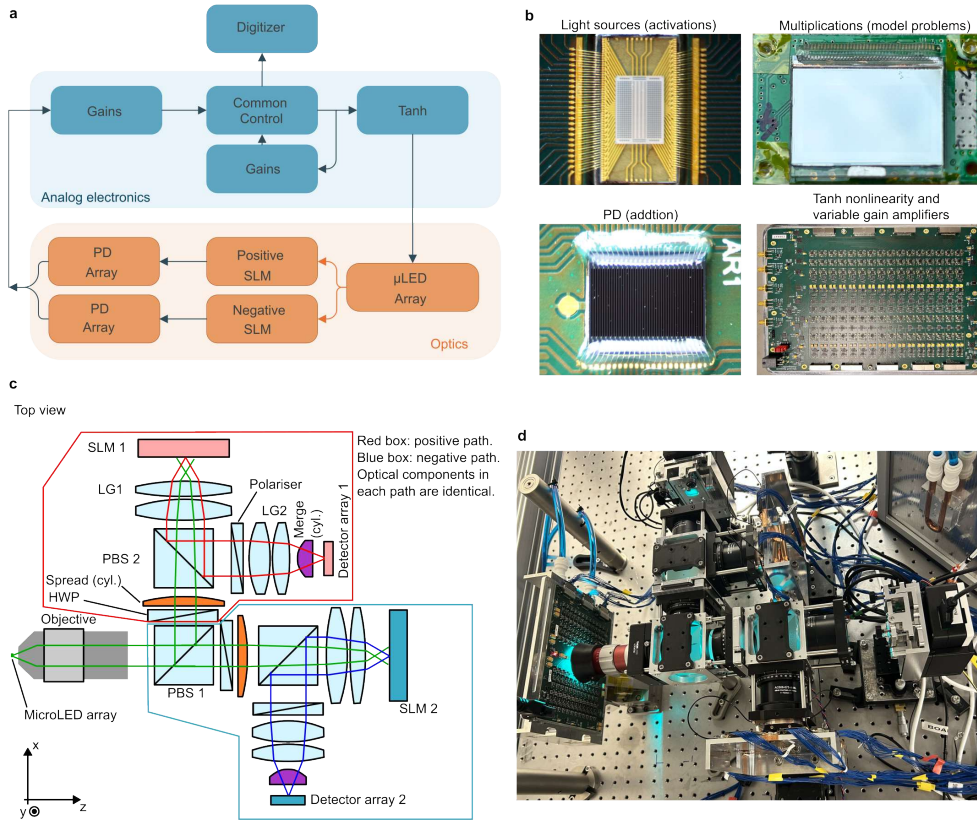
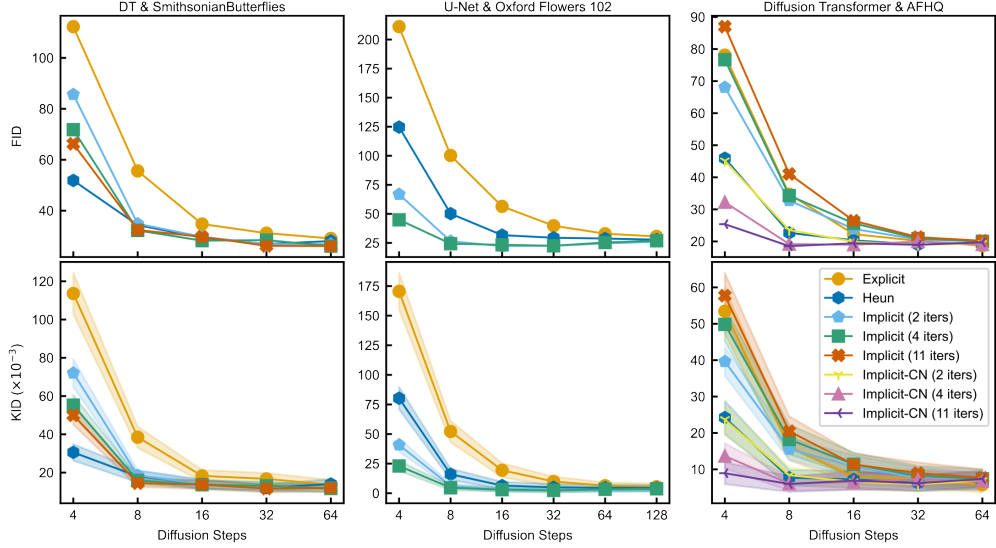


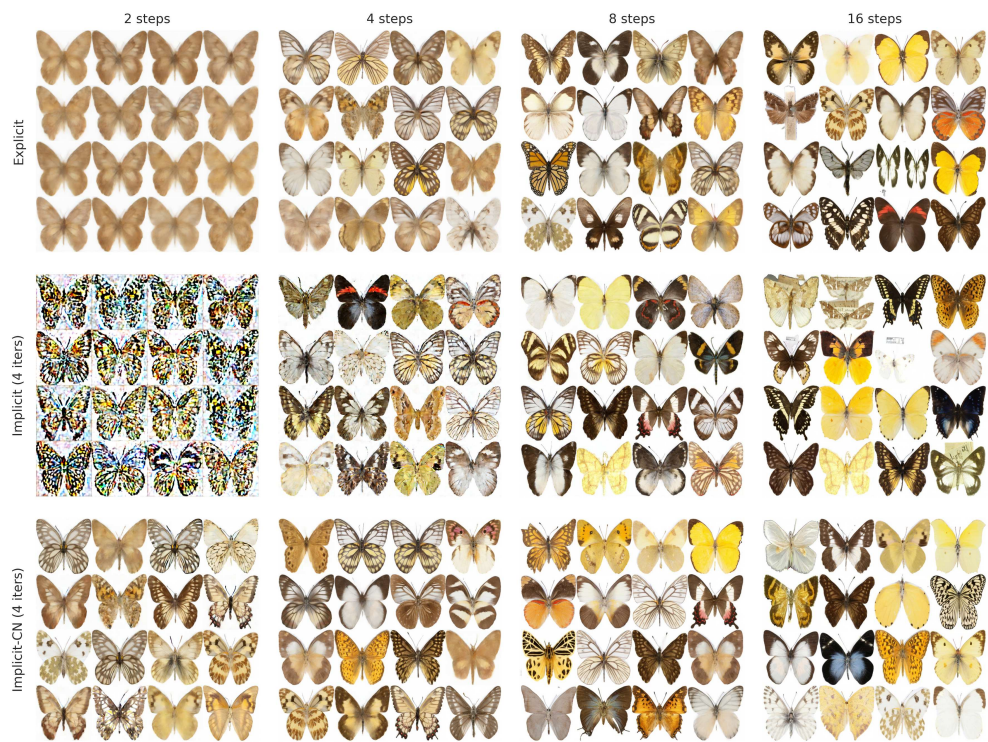
Fig. 4: Benchmarking ADMs at scale. **a**, FID versus number of diffusion steps for variants of implicit (this work) and explicit (standard) sampling from the same diffusion models for multiple datasets and model architectures. Implicit models use 4 fixed-point iterations per diffusion step. **b**, Qualitative comparison of sampling schemes at low step numbers. Butterflies (DT) and Flowers (U-Net) are generated with the standard Implicit method and AFHQ (DiT) shows the Implicit-CN variant. **c**, Trade-off study between sample quality and energy consumption: each curve represents the FID at 8, 16, 32 and 64 diffusion steps respectively, versus the normalized energy for different numbers of fixed-point iterations per diffusion step. **d**, Qualitative demonstration of Implicit-CN sampling with the state-of-the-art text-to-image diffusion model FLUX.1⁵⁵ using the prompt “a guitar”.



Extended Data Fig. 1: Hardware architecture of the analog computing platform and corresponding key components. **a**, Schematic diagram of the analog electronic (in blue) and the optical (in orange) sub-systems, implementing the fixed-point update rule (Eq.1). Arrows are shown to indicate the signals' flow. **b**, Key hardware components photos. The 4x48 microLED array (top left) is the light source and represents neural network activations. The spatial light modulator (SLM) (top right) stores neural network weights, and multiplies them with the incoming light. The 48 photodetector array (bottom left) adds and transfers optical signals into the analog electronic domain. The nonlinearity, programmable variable gains, and other computations are applied in analog electronics. **c**, Schematic of the optical subsystem, illustrating dual optical paths for positive and negative weights. The sources, collected by the objectives, are imaged onto the SLM using the combined 4F system of the objective and lens groups, respectively. **d**, Overhead view of the optical analog hardware.



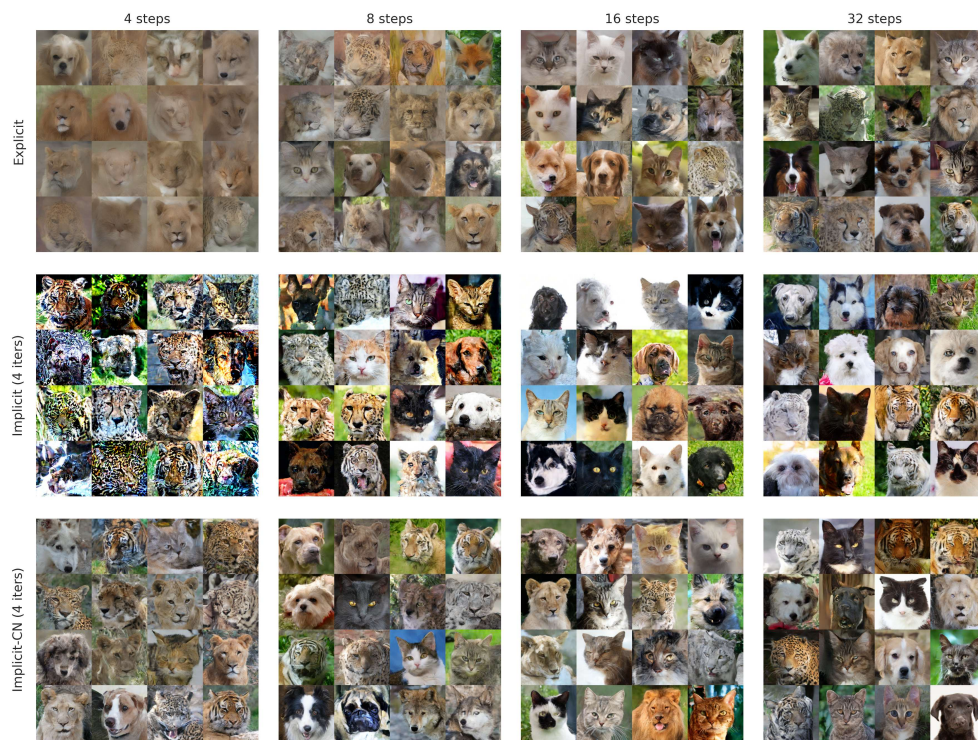
Extended Data Fig. 2: FID and Kernel Inception Distance (KID) for additional sampler configurations. For 8 or more diffusion steps 2 fixed-point iterations can be sufficient to reach convergence of the implicit solvers. Implicit sampling can yield the same (on Butterflies with ≥ 8 steps) or better performance (on Flowers) with the same number of neural network evaluations compared to a higher order Heun solver, which is a common best practice³³. Implicit-CN with 2 fixed-point iterations is mathematically identical to the Heun scheme and additional iterations can provide further improvements as shown for AFHQ. In all experiments we use a warm start of the fixed-point search, by initializing at the explicit Euler result, which requires one network evaluation and hence is counted as one fixed-point iteration.



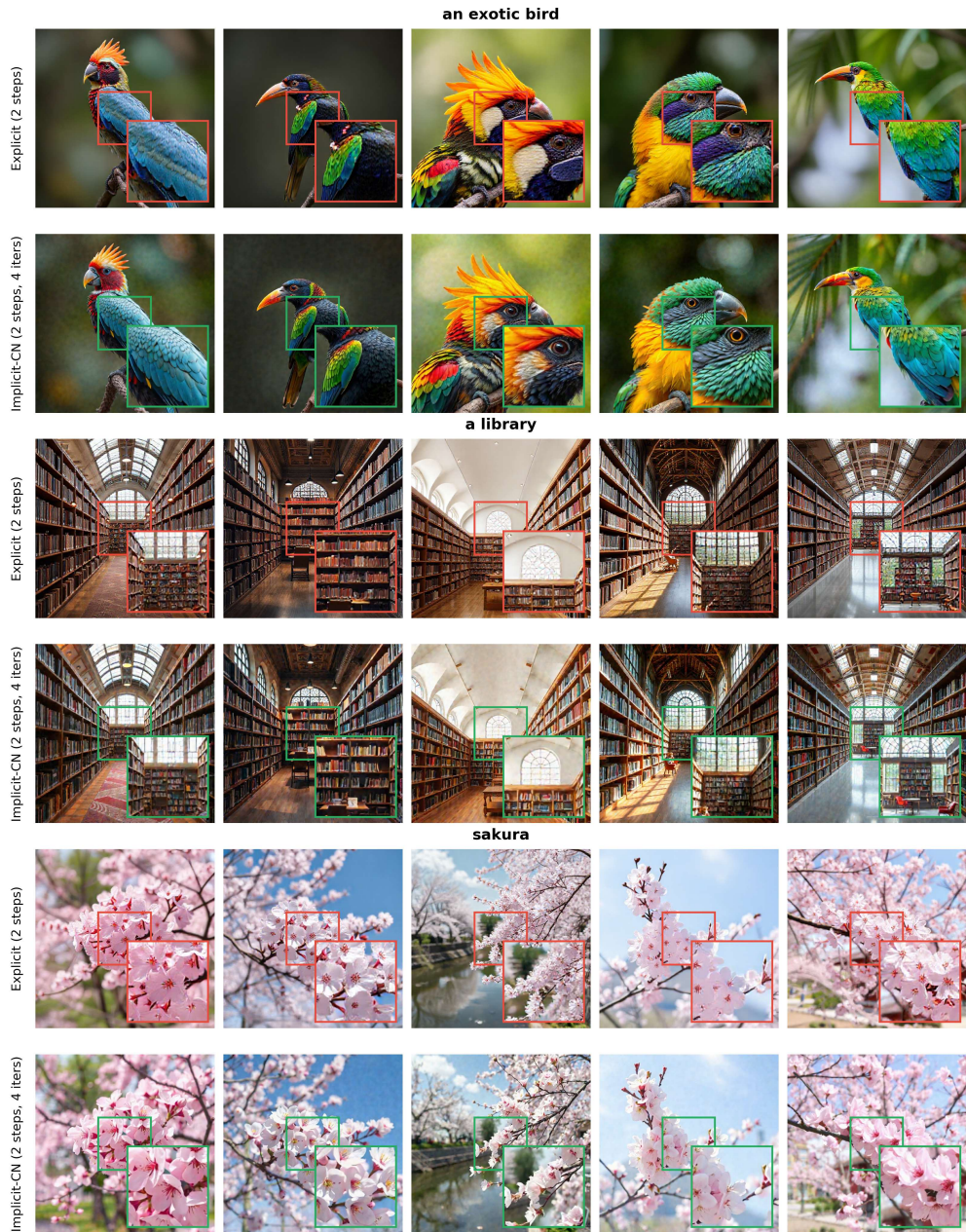
Extended Data Fig. 3: Additional samples for the three schemes (Explicit, Implicit and Implicit-CN) from the Digital Twin trained on the Smithsonian Butterflies dataset.



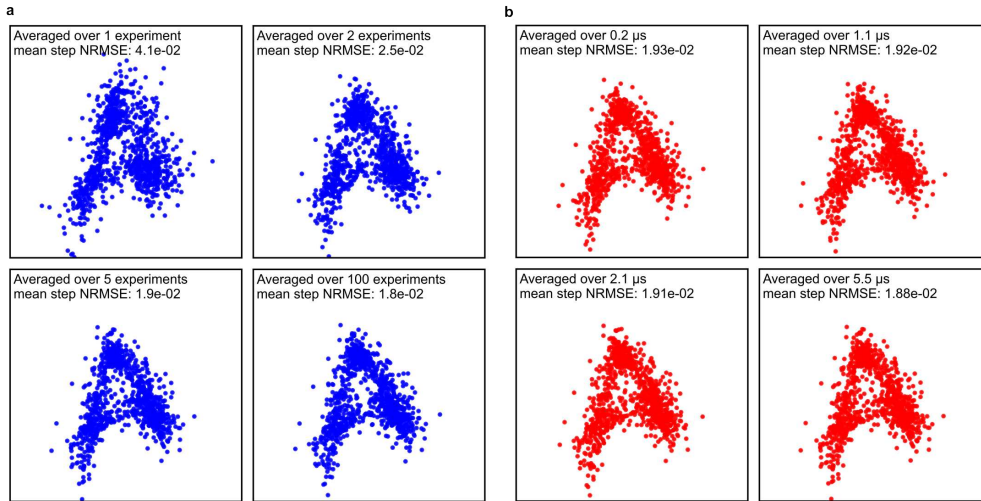
Extended Data Fig. 4: Additional samples for the three schemes (Explicit, Implicit and Implicit-CN) from the U-Net trained on Oxford-Flowers102.



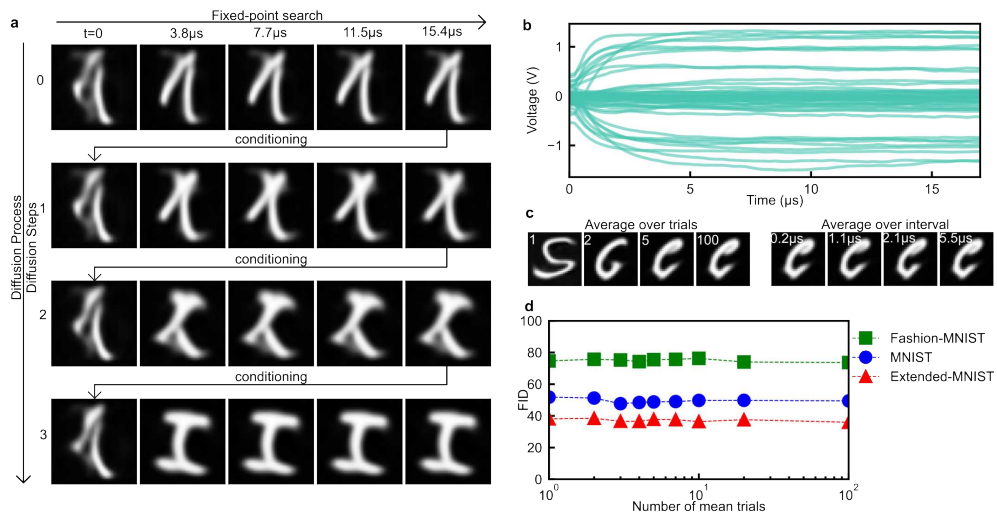
Extended Data Fig. 5: Additional samples for the three schemes (Explicit, Implicit and Implicit-CN) from the Diffusion Transformer (DiT) trained on AFHQ.



Extended Data Fig. 6: Additional samples for different prompts produced with the state-of-the-art text-to-image diffusion model FLUX.1^[55] using explicit Euler (default) and Implicit-CN sampling with 2 diffusion steps and 4 fixed-point iterations. The respective prompts are contained in the titles.



Extended Data Fig. 7: Additional noise details on the 2D experiments. **a**, Effect of averaging across repeated 2D experiments for the letter A, chosen as a representative example, for ADMs generated on AOC. Increasing the averaging beyond 5 yields negligible additional reductions in NRMSE. **b**, Effect of averaging within a fixed-point time window for the letter A for ADMs generated on AOC. Increasing the averaging window (up to 5.5 μ s) yields small but consistent reductions in NRMSE, consistent with SLM-dominated noise.



Extended Data Fig. 8: Inference dynamics for latent ADM hardware experiments on EMNIST, chosen as a representative example among the considered datasets. **a**, Evolution of latent ADM inference for EMNIST. Each row shows the progression of the fixed-point search within a diffusion step, followed by successive diffusion steps (rows). **b**, Example traces of the 48 analog voltages over time for a representative fixed-point within a diffusion step. **c**, Effect of averaging across repeated experiments and within a fixed-point time window. **d**, FID as a function of the number of mean trials for MNIST, FMNIST and EMNIST, showing negligible improvement beyond five averages across datasets.

577 References

- 578 [1] Thompson, N. C., Greenewald, K., Lee, K., Manso, G. F. *et al.* The computational
579 limits of deep learning. *arXiv preprint arXiv:2007.05558* **10**, 2 (2020).
- 580 [2] PILZ, K. F., MAHMOOD, Y. & HEIM, L. Ai's power requirements under
581 exponential growth (2025).
- 582 [3] Horowitz, M. *et al.* Scaling, power, and the future of cmos 7–pp (2005).
- 583 [4] Ahmed, S. R. *et al.* Universal photonic artificial intelligence acceleration. *Nature*
584 **640**, 368–374 (2025).
- 585 [5] Chen, Y. *et al.* All-analog photoelectronic chip for high-speed vision tasks. *Nature*
586 **623**, 48–57 (2023).
- 587 [6] Dong, B. *et al.* Partial coherence enhances parallelized photonic computing.
588 *Nature* **632**, 55–62 (2024).
- 589 [7] Hua, S. *et al.* An integrated large-scale photonic accelerator with ultralow latency.
590 *Nature* **640**, 361–367 (2025).
- 591 [8] Leroux, N. *et al.* Analog in-memory computing attention mechanism for fast
592 and energy-efficient large language models. *Nature Computational Science* 1–12
593 (2025).
- 594 [9] Zhang, Y. *et al.* Direct tensor processing with coherent light. *Nature Photonics*
595 1–7 (2025).
- 596 [10] Shen, Y. *et al.* Deep learning with coherent nanophotonic circuits. *Nature*
597 *photonics* **11**, 441–446 (2017).

- 598 [11] Oguz, I. *et al.* Optical diffusion models for image generation. *Advances in Neural*
599 *Information Processing Systems* **37**, 59150–59173 (2024).
- 600 [12] Chen, S., Li, Y., Wang, Y., Chen, H. & Ozcan, A. Optical generative models.
601 *Nature* **644**, 903–911 (2025).
- 602 [13] Cheng, Y. *et al.* Voltage-controlled magnetoelectric devices for neuromorphic
603 diffusion process. *Nature Communications* **16**, 5022 (2025).
- 604 [14] Yang, J. *et al.* Resistive memory-based neural differential equation solver for
605 score-based diffusion model. *arXiv preprint arXiv:2404.05648* (2024).
- 606 [15] Shastri, B. J. *et al.* Photonics for artificial intelligence and neuromorphic
607 computing. *Nature Photonics* **15**, 102–114 (2021).
- 608 [16] Feldmann, J., Youngblood, N., Wright, C. D., Bhaskaran, H. & Pernice, W. H.
609 All-optical spiking neurosynaptic networks with self-learning capabilities. *Nature*
610 **569**, 208–214 (2019).
- 611 [17] Wetzstein, G. *et al.* Inference in artificial intelligence with deep optics and
612 photonics. *Nature* **588**, 39–47 (2020).
- 613 [18] Hamerly, R., Bernstein, L., Sludds, A., Soljačić, M. & Englund, D. Large-scale
614 optical neural networks based on photoelectric multiplication. *Physical Review X*
615 **9**, 021032 (2019).
- 616 [19] Wang, T. *et al.* An optical neural network using less than 1 photon per
617 multiplication. *Nature Communications* **13**, 123 (2022).
- 618 [20] Kalinin, K. P. *et al.* Analog optical computer for ai inference and combinatorial
619 optimization. *Nature* 1–8 (2025).

- 620 [21] Sohl-Dickstein, J., Weiss, E., Maheswaranathan, N. & Ganguli, S. Deep unsuper-
621 vised learning using nonequilibrium thermodynamics. *International conference*
622 *on machine learning* 2256–2265 (2015).
- 623 [22] Ho, J., Jain, A. & Abbeel, P. Denoising diffusion probabilistic models. *Advances*
624 *in neural information processing systems* **33**, 6840–6851 (2020).
- 625 [23] Baldrige, J. *et al.* Imagen 3. *arXiv preprint arXiv:2408.07009* (2024).
- 626 [24] Polyak, A. *et al.* Movie gen: A cast of media foundation models. *arXiv preprint*
627 *arXiv:2410.13720* (2024).
- 628 [25] DeepMind, G. Veo 3. <https://deepmind.google/models/veo/> (2024).
- 629 [26] Blattmann, A. *et al.* Stable video diffusion: Scaling latent video diffusion models
630 to large datasets. *arXiv preprint arXiv:2311.15127* (2023).
- 631 [27] Lewis, S. *et al.* Scalable emulation of protein equilibrium ensembles with
632 generative deep learning. *Science* (2025).
- 633 [28] Zeni, C. *et al.* A generative model for inorganic materials design. *Nature* **639**,
634 624–632 (2025).
- 635 [29] Ivanov, A., Dryden, N., Ben-Nun, T., Li, S. & Hoefler, T. Data
636 movement is all you need: A case study on optimizing transform-
637 ers (2021). URL [https://proceedings.mlsys.org/paper_files/paper/2021/file/
638 bc86e95606a6392f51f95a8de106728d-Paper.pdf](https://proceedings.mlsys.org/paper_files/paper/2021/file/bc86e95606a6392f51f95a8de106728d-Paper.pdf).
- 639 [30] Song, Y., Dhariwal, P., Chen, M. & Sutskever, I. Consistency models (2023).
- 640 [31] Liu, X., Gong, C. & Liu, Q. Flow straight and fast: Learning to generate and
641 transfer data with rectified flow. *arXiv preprint arXiv:2209.03003* (2022).

- 642 [32] Song, J., Meng, C. & Ermon, S. Denoising diffusion implicit models. *arXiv*
643 *preprint arXiv:2010.02502* (2020).
- 644 [33] Karras, T., Aittala, M., Aila, T. & Laine, S. Elucidating the design space of
645 diffusion-based generative models. *Advances in neural information processing*
646 *systems* **35**, 26565–26577 (2022).
- 647 [34] Lu, C. *et al.* Dpm-solver: A fast ode solver for diffusion probabilistic model
648 sampling in around 10 steps. *Advances in neural information processing systems*
649 **35**, 5775–5787 (2022).
- 650 [35] Song, Y. *et al.* Score-based generative modeling through stochastic differential
651 equations. *arXiv preprint arXiv:2011.13456* (2020).
- 652 [36] Bai, S., Kolter, J. Z. & Koltun, V. Deep equilibrium models. *Advances in neural*
653 *information processing systems* **32** (2019).
- 654 [37] Geng, Z., Pokle, A. & Kolter, J. Z. One-step diffusion distillation via deep
655 equilibrium models. *Advances in Neural Information Processing Systems* **36**,
656 41914–41931 (2023).
- 657 [38] Pokle, A., Geng, Z. & Kolter, J. Z. Deep equilibrium approaches to diffusion
658 models. *Advances in Neural Information Processing Systems* **35**, 37975–37990
659 (2022).
- 660 [39] Bai, X. & Melas-Kyriazi, L. Fixed point diffusion models (2024).
- 661 [40] Syed, M., Kalinin, K. & Berloff, N. Beyond digital: Harnessing analog hardware
662 for machine learning (2023).
- 663 [41] Hopfield, J. J. Neurons with graded response have collective computational prop-
664 erties like those of two-state neurons. *Proceedings of the national academy of*

- 665 *sciences* **81**, 3088–3092 (1984).
- 666 [42] Ambrogio, S. *et al.* An analog-ai chip for energy-efficient speech recognition and
667 transcription (2023).
- 668 [43] Zoppo, G., Marrone, F. & Corinto, F. Equilibrium propagation for memristor-
669 based recurrent neural networks (2020).
- 670 [44] Al-Kayed, N. *et al.* Programmable 200 gops hopfield-inspired photonic ising
671 machine. *Nature* **648**, 576–584 (2025).
- 672 [45] Momeni, A. *et al.* Training of physical neural networks. *Nature* **645**, 53–61 (2025).
- 673 [46] Wright, L. G. *et al.* Deep physical neural networks trained with backpropagation.
674 *Nature* **601**, 549–555 (2022).
- 675 [47] Rasch, M. J. *et al.* Hardware-aware training for large-scale and diverse deep learn-
676 ing inference workloads using in-memory computing-based accelerators. *Nature*
677 *communications* **14**, 5282 (2023).
- 678 [48] Rombach, R., Blattmann, A., Lorenz, D., Esser, P. & Ommer, B. High-resolution
679 image synthesis with latent diffusion models 10684–10695 (2022).
- 680 [49] Ronneberger, O., Fischer, P. & Brox, T. U-net: Convolutional networks for
681 biomedical image segmentation (2015).
- 682 [50] Peebles, W. & Xie, S. Scalable diffusion models with transformers (2023).
- 683 [51] Lipman, Y., Chen, R. T., Ben-Hamu, H., Nickel, M. & Le, M. Flow matching for
684 generative modeling. *arXiv preprint arXiv:2210.02747* (2022).
- 685 [52] Albergo, M. S. & Vanden-Eijnden, E. Building normalizing flows with stochastic
686 interpolants. *arXiv preprint arXiv:2209.15571* (2022).

- 687 [53] Jolicoeur-Martineau, A., Li, K., Piché-Taillefer, R., Kachman, T. & Mitliagkas,
688 I. Gotta go fast when generating data with score-based models. *arXiv preprint*
689 *arXiv:2105.14080* (2021).
- 690 [54] Shebanow, M., Finkelstein, H. & Bowen, P. Revolutionizing the ai
691 data center: Delivering 4,000 peta operations per second at 1% of the
692 power using optical systolic arrays. White Paper, Neurophos Inc. (2025).
693 URL [https://11549dc2-0dc9-4363-9a47-d257d2a497cc.filesusr.com/ugd/79aa3c_](https://11549dc2-0dc9-4363-9a47-d257d2a497cc.filesusr.com/ugd/79aa3c_85cd1a5477d44f4fa0e423c03ad0746c.pdf)
694 [85cd1a5477d44f4fa0e423c03ad0746c.pdf](https://11549dc2-0dc9-4363-9a47-d257d2a497cc.filesusr.com/ugd/79aa3c_85cd1a5477d44f4fa0e423c03ad0746c.pdf). Accessed: 2026-01-28.
- 695 [55] Labs, B. F. *et al.* Flux. 1 kontekst: Flow matching for in-context image generation
696 and editing in latent space. *arXiv preprint arXiv:2506.15742* (2025).
- 697 [56] Gadhikar, A., Grazi, R., Hensman, J. *et al.* Optrot: Mitigating weight out-
698 liers via data-free rotations for post-training quantization (2024). URL [https:](https://openreview.net/forum?id=b5649238190ee2a0c8b8d7c25154a090163b6b4b)
699 [//openreview.net/forum?id=b5649238190ee2a0c8b8d7c25154a090163b6b4b](https://openreview.net/forum?id=b5649238190ee2a0c8b8d7c25154a090163b6b4b).
- 700 [57] Farhat, N. H., Psaltis, D., Prata, A. & Paek, E. Optical implementation of the
701 hopfield model. *Applied optics* **24**, 1469–1475 (1985).
- 702 [58] Glorot, X. & Bengio, Y. Understanding the difficulty of training deep feedforward
703 neural networks (2010).
- 704 [59] Tolstikhin, I. O. *et al.* Mlp-mixer: An all-mlp architecture for vision. *Advances*
705 *in neural information processing systems* **34**, 24261–24272 (2021).

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [Supplementary.pdf](#)