

---

# Optical Transformers

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 The rapidly increasing size of deep-learning models has caused renewed and grow-  
2 ing interest in alternatives to digital computers to dramatically reduce the energy  
3 cost of running state-of-the-art neural networks. Optical matrix-vector multipliers  
4 are best suited to performing computations with very large operands, which leads  
5 us to hypothesize that large Transformer models might achieve asymptotic energy  
6 advantages with optics over running digitally. To test this idea, we performed  
7 small-scale optical experiments with a prototype accelerator to demonstrate that  
8 Transformer operations can run on optical hardware despite noise and errors. Using  
9 experiment-calibrated simulations of our hardware, we studied the behavior of  
10 running Transformers optically, identifying scaling laws for model performance  
11 with respect to optical energy usage and estimating total system power consump-  
12 tion. We found that the optical energy per multiply-accumulate (MAC) scales as  
13  $\frac{1}{d}$  where  $d$  is the Transformer width, an asymptotic advantage over digital sys-  
14 tems. Should well-engineered, large-scale optical hardware be developed, it might  
15 achieve a  $100\times$  energy-efficiency advantage for running some of the largest current  
16 Transformer models, and if both the models and the optical hardware are scaled  
17 to the quadrillion-parameter regime, optical computers could have a  $> 8,000\times$   
18 energy-efficiency advantage over state-of-the-art digital-electronic processors (300  
19 fJ/MAC). We discussed how these results motivate and inform the construction of  
20 future optical accelerators and optics-amenable deep-learning approaches. With  
21 assumptions about future improvements to electronics and Transformer quantiza-  
22 tion techniques ( $5\times$  cheaper memory access, double the digital-analog conversion  
23 efficiency, and 4-bit precision), we estimated that optical computers' advantage  
24 against these digital processors could grow to  $> 100,000\times$ .

## 25 1 Introduction

26 Deep learning models' exponentially increasing scale is both a key driver in advancing the state-of-  
27 the-art and a cause of growing concern about their energy usage, speed, and practicality. This has led  
28 to the development of hardware accelerators and model training/compression/design techniques for  
29 efficient and fast inference on them.

30 While digital-electronic accelerators [47, 16, 8, 1, 17] can improve performance by some constant  
31 factor, alternative analog computing platforms using optics have been proposed as a new paradigm  
32 for better scalability [49, 7, 62, 41, 56, 24, 51]. Ideally, the scaling is asymptotically better than  
33 digital systems in energy per MAC [18, 61, 53, 41]. But these optical neural networks (ONNs) have  
34 additional complexities and limitations of their own such as low precision, noise, and analog/digital  
35 data conversion overheads which depend on the access patterns of the model running (Figure 1).  
36 Thus, advantageously accelerating any neural network architecture with ONNs is hard. Here, we  
37 hope to answer whether Transformers' efficient data-access patterns (wide layers, parallel/batched

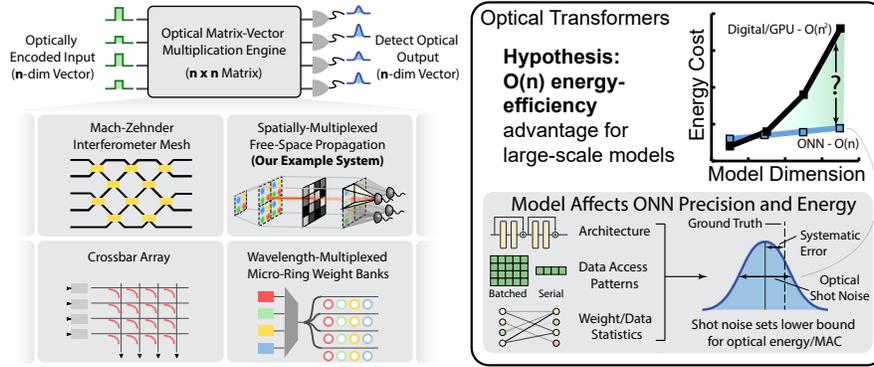


Figure 1: **Can Transformers Benefit From Running on Optical Hardware?** Optical Neural Networks (ONNs) have been proposed as an alternative computing platform that can achieve asymptotic energy-efficiency advantages over digital computers running neural networks. This is not a guarantee; their behavior is affected by model architecture, statistics, and resilience to the noise/imprecision of analog hardware. Thus, while there are many implementations of general-purpose optical matrix accelerators (such as those depicted in the inset), there are still model-dependent challenges/tradeoffs in realizing their purported advantages. We seek here to answer the question of how much today’s enormous Transformer models can benefit from this technology, if at all. Our hypothesis is that Transformers’ architecture and unique behaviors allow for ONN-enabled benefits that scale.

38 token processing, etc.), trends in methods for scaling them, and sufficient effort to train them for  
 39 ONNs afford them the asymptotic energy-efficiency advantages of running optically.

40 Here we demonstrate how the popular Transformer architecture is able to run on ONN systems,  
 41 and estimate the potential benefits of doing so. To first verify that Transformers may run on these  
 42 systems despite their imprecision, we sampled operations from a Transformer and ran them on a real  
 43 spatial light modulator (SLM) based experimental system, and used the results to create a calibrated  
 44 simulation of the optical hardware, with the systematic error, noise, and imprecision of weights/inputs  
 45 we observed. Transformers running on the simulated hardware could perform nearly as well as those  
 46 running digitally, and could be far more efficient. We summarize our key contributions as follows:

- 47 • We demonstrated linear Transformer operations (the bulk of a Transformer’s computation)  
 48 running with sufficient accuracy on real optical hardware and in a matching simulation,  
 49 despite errors and noise.
- 50 • Via simulation, we established scaling laws for optical Transformer performance versus  
 51 optical energy usage, and optical energy usage versus model size.
- 52 • Based on our simulations and experiments we estimated an orders-of-magnitude energy  
 53 consumption advantage of full ONN accelerators versus state-of-the-art GPUs.
- 54 • We discussed Transformers’ suitability for optical acceleration, and more generally how  
 55 specific elements of DNN architecture affect the function of ONN systems running them.
- 56 • We identified the hardware and systems design challenges that future work on building ONN  
 57 accelerators should target.

58 While our experiments and simulations were based on specific hardware as a representative example,  
 59 our scope here is more general. We are interested in understanding how uniquely optical energy  
 60 scaling and noise relate to Transformer performance and architecture. As such nearly all our findings  
 61 apply broadly to linear optical processors (and hopefully future ones), irrespective of their underlying  
 62 hardware implementation details.

## 63 2 Background and Related Work

### 64 2.1 Transformer Models

65 Transformers are models for processing sequential data based on multi-head attention. Transformers  
 66 consist of two-layer feed-forward blocks and multi-head attention (Figure 2) operations. Multi-

67 head attention computes relationships between sequence elements by deriving query, key, and  
 68 value sequences  $Q, K, V$  and computing dot products with a softmax nonlinearity in-between [60].  
 69 Transformers also leverage modern design elements such as additive residual skip connections [20]  
 70 and normalization layers [3]. A defining feature of Transformers is that entire sequences may be  
 71 processed in matrix-matrix products in parallel (instead of one token/input at a time).

## 72 2.2 Large-Scale Deep Learning

73 In the past few years, it has been found in particular that Transformer [60] architectures significantly  
 74 improve when sized up to billions or even trillions of parameters [6, 28, 10, 22, 59, 66], causing an  
 75 exponential growth of deep learning compute usage [48, 50]. These large-scale Transformers achieve  
 76 ever more impressive results in not only natural language processing, but also in other domains such  
 77 as computer vision [14, 36], graphs [30], and in multi-modal settings [27, 26, 44, 45, 65, 46], making  
 78 them a popular but expensive solution for many tasks—digital hardware’s energy efficiency (ie.  
 79 per-flop or per-inference cost) has not kept up with the growing FLOP requirements of state-of-the-art  
 80 deep learning models [50]. They also have transfer learning capabilities [42, 13, 43, 6, 37, 14],  
 81 allowing them to easily generalize to specific tasks, in some cases in a zero-shot setting where no  
 82 further training is necessary [6, 45, 33].

## 83 2.3 Optical Accelerators

84 Researchers have explored a wide variety of controllable optical systems which manipulate different  
 85 types of optical modes to effectively implement arbitrary matrix-vector multiplications, vector-vector  
 86 dot products [52, 2, 18, 55, 4, 61, 19, 39, 57], or convolutions [63, 15, 40, 64]. In this work, we adopt  
 87 the free-space multiplier [61, 55, 19] (Figure 2, top left) to demonstrate Transformer operations in  
 88 optical experiments and for our simulations. We selected this system because it has many of the same  
 89 behaviors as other ONN implementations, and aim to draw conclusions that could generally be useful  
 90 for those working with other ONN designs. Many ONN systems, including ours, share the following  
 91 typical traits:

92 **Device Imprecision and Optical Shot Noise** Optical systems are subject to errors in both the  
 93 actual hardware and from photon detection. Detection of optical intensity in particular is subject to a  
 94 phenomenon known as *shot noise* where the detected value is Poisson distributed: given vectors  $x$   
 95 and  $w$ , with the elements of  $x$  encoded as optical intensity, the output  $Y$  is distributed as:

$$Y \sim \text{Poisson}(w \cdot x) \tag{1}$$

96 For other encoding schemes such as amplitude or phase encoding, equation 1 should be modified, but  
 97 the detection is still subject to shot noise.

98 **Efficient Photon Usage** Shot noise, and therefore an optical dot product’s signal-to-noise ratio  
 99 (SNR, which serves as an effective bit precision) is related to the mean number of photons at the  
 100 *output*. The efficiency of photon usage can therefore grow with increasing multiply-accumulate  
 101 operations (MACs): the SNR for the product  $w \cdot x$  is

$$\text{SNR}(Y) = \frac{\text{E}[Y]}{\sqrt{\text{Var}[Y]}} = \sqrt{w \cdot x} = \sqrt{\text{E}[Y]}, \tag{2}$$

102 which explains this behavior; if the desired output precision does not change, constant photons are  
 103 required regardless of dot product size. Work on ONNs has studied this behavior in a variety of  
 104 scenarios [18, 41, 61, 53]. This efficient scaling is not a guarantee—the required number of photons  
 105 may be influenced by a model architecture’s activation/weight distributions, encoding schemes,  
 106 precision requirements, etc.

107 **Optical Neural Network Energy Costs** The energy cost of optical neural networks is broken down  
 108 into the optical costs of performing MACs and the electrical costs of loading/detecting data, which  
 109 are usually dominant. Consider a product between two matrices,  $A \in \mathbb{R}^{n \times d}$ ,  $B \in \mathbb{R}^{d \times k}$ . Such a  
 110 product results in loading (detecting)  $nd + dk$  ( $nk$ ) scalars, and performing  $ndk$  MACs. If the energy

111 to electrically load (detect) a scalar is  $E_{\text{load}}$  ( $E_{\text{det}}$ ), and to perform a MAC optically is  $E_{\text{optical}}$ , then  
112 the total energy is:

$$E = (nd + dk)E_{\text{load}} + nkE_{\text{det}} + ndkE_{\text{optical}} \quad (3)$$

113 This illustrates how ONNs may have asymptotic energy advantages over digital computers. Notice  
114 that regardless of the number of reuses, all data is only loaded once in Equation 3. This is because  
115 copying a vector’s data and transporting it is free optically. Meanwhile,  $E_{\text{optical}}$  ideally scales as  $1/d$ .  
116 These properties make energy cost disproportional to the number of MACs,  $ndk$ . In other words,  
117  $\frac{E_{\text{digital}}}{E_{\text{ONN}}} \sim \min(n, d)$ .

118 **Streaming Weights Versus Weights-In-Place** There are two approaches for loading  
119 weights. *Weights-in-place* schemes involve loading them once, and re-using them for many inputs.  
120 Alternatively, systems can employ *streaming weights* where at every computation the required weight  
121 matrix is loaded. Our experimental system is a weights-in-place scheme. For weights-in-place  
122 operations, the energy advantage scales as just  $\frac{E_{\text{digital}}}{E_{\text{ONN}}} \sim d$ .

## 123 2.4 Previous Optical Neural Network Architectures

124 Previous work has considered deep learning models such as MLPs and convolutional networks  
125 on benchmark tasks like MNIST [40, 61], and simulations of larger convolutional models such as  
126 AlexNet [32] on more difficult datasets such as ImageNet [18]. This begs the question of how well  
127 newer, larger models perform on optical systems.

## 128 2.5 Scalable Compression and Quantization of Large Language Models (LLMs)

129 Optical hardware’s low precision raises the question of whether scaled-up models could be quantized  
130 sufficiently to run. Thankfully, continual research in LLM compression has progressively shown that  
131 larger models do not have increasing precision requirements. For example, [34] found that larger  
132 Transformers can be compressed more easily, to the degree that it is more worthwhile to train large  
133 ones and compress them over training smaller ones of the target size. Furthermore, [5] and [12]  
134 demonstrated running Transformers at scale with int8 precision, and the recent work of [11] proposes  
135 that 4-bit is optimal for nearly all model scales, except for the largest tested (175B parameters) where  
136 3-bit was sometimes found to work better.

# 137 3 Optical Transformers

138 We designed models that are intentionally similar to other Transformers, with the goal of simulating  
139 their behavior (informed by some experimental measurements) and energy consumption on optical  
140 hardware. A summary of our approach and model is in Figure 2.

## 141 3.1 Architecture and Task

142 We created optical Transformer models with a GPT2-like [43] architecture that replaces the GELU  
143 [21] activation with ReLU6, which is known to improve low-precision model performance [31, 23, 29].  
144 For language modelling, we used the raw Wikitext-103 dataset [38]. The models we simulated have  
145 12 layers (consisting of multi-head attention and feed-forward blocks), operate on a context length  
146 of 1024 tokens, use 12 attention heads, and have embedding dimension  $d$  varying from 192 to 1536.  
147 The full details of the training technique, architecture, and hyperparameters are in Appendix A.

## 148 3.2 Transformer Computations on Optical Hardware

149 We ran experiments using a real Transformer’s (we used the base-sized model with  $d = 768$ ) weights  
150 in order to characterize the behavior of an ONN system. We adopted as a representative example of  
151 an optical accelerator a spatial light modulator (SLM) based system which computes vector-vector  
152 dot products [61]. Vectors are encoded on a display, and copies are shone through the SLM which  
153 has varying transmission corresponding to some data (ie. a weight matrix). The outputs of this

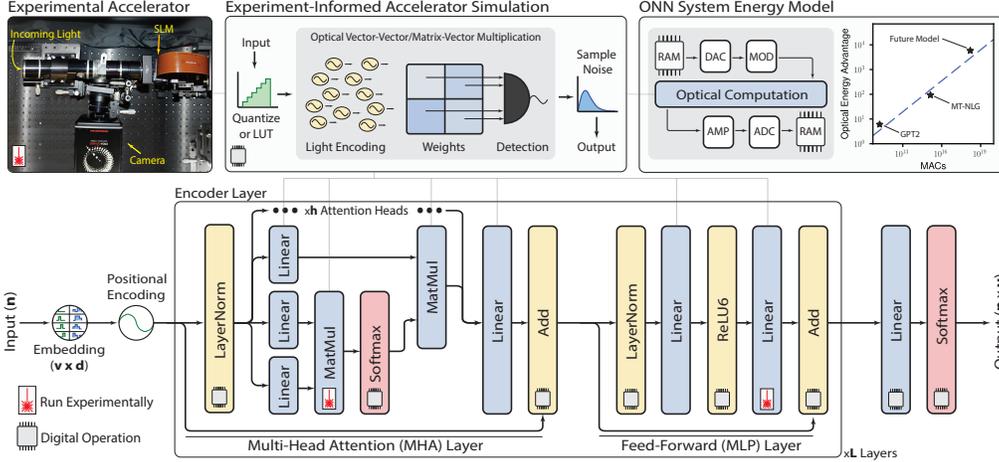


Figure 2: **Optical Transformer evaluation: prototype hardware; simulator model; Transformer architecture.** Bottom: typical Transformer architecture, but with ReLU6 activation. Top Left: experimental spatial light modulator (SLM)-based accelerator setup. From some layers—marked with a laser icon—we sampled dot products to run on real hardware. Top Middle: Linear operations, in light blue, run on a simulated accelerator with noise/error. Lookup tables (LUT) allow simulation using our setup’s supported weight/activation values. Top right: our model of energy consumption for optical accelerators, based on assumptions and results from our experiment/simulations. The model accelerator system consists of random-access memory (RAM), a analog/digital conversion (DAC/ADC), light modulation (MOD), amplification (AMP).

154 operation—element-wise products—are collected at detectors as the resultant dot products (Figure 2,  
 155 top left). We collected lookup tables (LUTs)—mappings of the available discrete levels in both the  
 156 display and SLM devices—and used them to train a “LUT-aware” optical Transformer model to run  
 157 on the setup. We then collected calibration curves, mappings from the detected output light intensity  
 158 to the actual neuron floating-point values. To do this, we ran many random dot products on the  
 159 hardware and collected pairs of detected values and digitally-computed ground-truth values. We then  
 160 fit the relationship linearly. We used high photon counts to eliminate shot noise, so deviation from  
 161 the linear fit was considered the hardware’s *systematic error*. Full details of experimental procedures  
 162 and calibration are in Appendix B.

### 163 3.3 Simulation of Optical Hardware

164 Informed by our experiments, we  
 165 constructed a simulation of the  
 166 optical hardware. By simulat-  
 167 ing the hardware behavior di-  
 168 rectly we model how any arbi-  
 169 trary operation would behave if  
 170 run on the physical setup. This  
 171 allows us to avoid the computa-  
 172 tionally demanding task of simul-  
 173 ating much larger Transformers  
 174 to verify that our simulation  
 175 method works. We aimed to em-  
 176 ulate the noise, error, and preci-  
 177 sion that we observed in order to understand how well full Transformers would perform when running  
 178 on optical hardware. The configurations for different scenarios are summarized in Table 1. We also  
 179 evaluated the digital, 8-bit-QAT-trained model for comparison purposes.

Table 1: Summary of simulation configurations for different evaluation and training scenarios. For simulating optical hardware we included all behaviors. For determining optical resource scaling, we focused on shot noise, and ran a plain 8-bit model for comparison.

Setting	Op.	Shot Noise	Sys. Err.	LUT	4-Pass
Hardware Simulation	QAT	✗	✗	✓	✗
	Eval	✓	✓	✓	✓
Optical Scaling Simulation	QAT	✗	✗	✗	✗
	Eval	✓	✗	✗	✓
	Int8	✗	✗	✗	✗

180 **Hybrid Scheme** Pure optical systems cannot easily compute activation or normalization functions.  
 181 Thus we assumed LayerNorm, ReLU activations, and residual skip connections are performed digitally  
 182 at full precision. Thankfully, even in smaller models, linear computations are the overwhelming  
 183 majority (Section 4.3).

184 **Non-Negative Weights and Inputs (“4-Pass” Multiplication)** An important limitation is that our  
 185 display and SLM only support non-negative values. The constraint of having all-positive data is  
 186 present in many but not all optical neural network systems. We worked around this by decomposing  
 187 products into sums/differences of products with non-negative operands. Consider a product between  
 188 matrices  $W$  and  $X$ . If we let  $W_+$  ( $X_+$ ) and  $W_-$  ( $X_-$ ) be matrices with only the positive and negative  
 189 elements of  $W$  ( $X$ ) respectively, then:

$$WX = W_+X_+ - |W_-|X_+ - W_+|X_-| + W_-X_- \quad (4)$$

190 **Data Scaling** On the real system, we define a maximum activation/weight value as 1.0 and minimum  
 191 as 0.0. To simulate operation, the inputs and weights of every simulated NN layer are scaled to this  
 192 range, and then rescaled back afterwards.

193 **Device Quantization** Real hardware may only have certain number of representable levels. To  
 194 emulate this behavior, we fine-tuned pretrained models using quantization-aware training [25](QAT)  
 195 and applied the following in simulation (hyperparameters in Appendix A):

- 196 • For optics-simulated layers, we emulated quantization to int8 (256 levels). Then, instead of  
 197 dequantizing, we used the integer values directly as indices into the LUTs that we gathered  
 198 from experiment.
- 199 • We also quantized weights, but with the SLM LUT. We clamped smaller values to 0.02 in the  
 200 simulation, as our SLM does not have a high extinction ratio, and the smallest transmission  
 201 is 0.02.
- 202 • Accumulation can be high precision, but we used int8 quantization for outputs, since  
 203 analog-digital conversion (ADC) is expensive in practice.
- 204 • We used both deterministic and stochastic rounding when quantizing, with similar results.

205 **Systematic Errors** Issues like cross-talk, misalignment, defects in ONNs give rise to systematic  
 206 errors. We simulated such a constraint by adding Gaussian noise to simulated model outputs  
 207 (Figure 2), scaled relative to the mean sizes of the outputs, as this was the noise behavior we observed  
 208 experimentally (it is related to the rescaling of data between 0 and 1).

209 **Optical Encoding and Shot Noise** We modeled optical encoding by subjecting layer outputs  
 210 to simulated shot noise (Figure 2), which differs from the systematic error model. Outputs were  
 211 scaled by a number such that the average photon number per feature (photons/MAC) was some  
 212 target value. Each of these features was used as the mean of a Poisson distribution, which we  
 213 sampled. These outputs were then scaled back down to represent neuron values. In the simulations  
 214 for optical scaling we used vanilla 8-bit QAT (no LUTs or systematic error, which can overwhelm  
 215 shot noise) to cleanly demonstrate the optical scaling properties—which are model-dependent and  
 216 not hardware-dependent—of Transformers.

## 217 4 Results

### 218 4.1 Transformer Error Tolerance and Hardware-Simulation Accuracy

219 We determined experimentally that Transformer operations are able to run on real hardware without  
 220 severely degraded performance from systematic errors. The bottom four panels of Figure 3 are  
 221 histograms of the experimental differences from correct values. The simulated noise distributions  
 222 (dotted lines) match well with the experimental data, which confirms that they are an accurate  
 223 representation of the real systematic error behavior. Figure 3 (top) is a map of the performance of the  
 224 simulated model over different configurations of the mean-relative (in percent) noise at every layer of  
 225 feed-forward and attention blocks. The model performs well with significant noise (experimental  
 226 noise levels marked with stars), within 1 perplexity from noise-free performance unless the noise is  
 227 very high. These results show that our digital model of the system is a plausible approximation of  
 228 how a real one might behave.

229 While 8-bit precision was used for  
 230 QAT, the optical Transformer can per-  
 231 form inference at lower precision, as  
 232 implied by its error tolerance. To  
 233 study this further we conducted a sim-  
 234 ple ablation on the input and output  
 235 precisions used at inference, on the 8-  
 236 bit-QAT base-sized model with LUT  
 237 in Appendix C.

## 238 4.2 Optical Scaling Laws

239 Optical Transformers achieve language  
 240 modelling performance close to their  
 241 digital counterparts’ when shot-noise-  
 242 limited at modest photon budgets.  
 243 The perplexities on the Wikitext-103  
 244 validation set of various optical Trans-  
 245 former models simulated with different  
 246 total photon usage (amount used for  
 247 input data) are shown in Figure 4 (left).  
 248 The curves illustrate a tradeoff: larger  
 249 models need larger photon totals to  
 250 function well, and there are different  
 251 optimal model choices based on the  
 252 photon budget. We define photons/MAC  
 253 as the total photon budget (amount at  
 254 input) divided by total MACs. The  
 255 percentage difference from the perfor-  
 256 mance at 10K photons/MAC (Figure 4,  
 257 middle)—chosen to represent an ideal  
 258 high-precision scenario—is roughly  
 259 power-law scaled in photons/MAC for  
 260 all models with truncation near 10K;  
 261 better performance can be had with  
 262 more photons, but with diminishing  
 263 returns, and the performance matches  
 264 or exceeds that of the 8-bit digital  
 265 models’ when the photon budget is  
 266 not too low ( $\sim 10^2$ ).

268 The models use fewer photons/MAC  
 269 as they scale, achieving the theoretical  
 270 efficient scaling where the total per-  
 271 dot-product photons needed is constant.  
 272 To study how photon usage scales, we  
 273 determined how many photons it takes  
 274 to reach the performance of 8-bit digital  
 275 models. These values, in Figure 4 (right),  
 276 decrease nearly as  $\frac{1}{d}$ —the total photons  
 277 needed per dot product is constant (bottom  
 278 dashed line). The Transformer archi-  
 279 tecture clearly takes advantage of effi-  
 280 cient optical scaling with larger model  
 281 sizes. In fact, smaller per-dot-product  
 282 totals are required for the largest model,  
 283 suggesting that larger Transformers  
 284 may require less output precision. This  
 285 is consistent with other work which  
 286 found that precision requirements are  
 287 constant or reduced with scale [34].  
 288 Meanwhile, the already low photon  
 289 usage of the largest model suggests that  
 290 models larger than our simulations  
 291 (>10B parameters) may use <1  
 292 photon/MAC. This sub-photon operation  
 293 works in optical systems [61, 53] and  
 294 is in essence no different at all from  
 295 operation at higher photon counts (since  
 296 the number summed at detection is still  
 297 high).

281 These empirical scaling results are tied  
 282 to our specific configurations and train-  
 283 ing strategies. Depending on the scales  
 284 and dynamic ranges of inputs and weights,  
 285 different amounts of photons may be  
 286 transmitted to the output; the statistics  
 287 of a model affect its efficiency. In  
 288 Appendix H we explore a different  
 289 scheme, but the effects of different  
 290 methods remains an interesting topic  
 291 for future work.

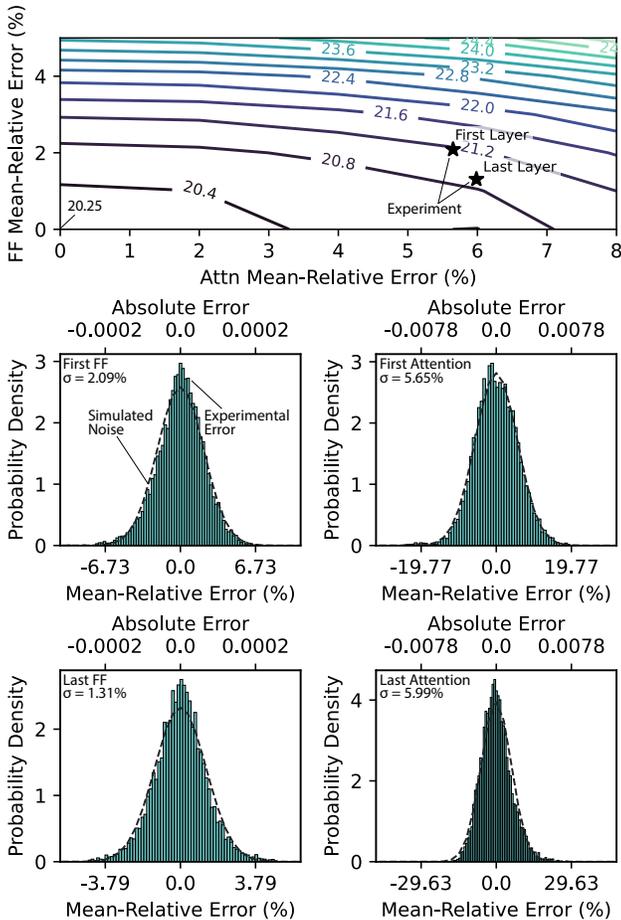


Figure 3: **Comparison of experimental and simulated noise models and simulated Optical Transformer noise tolerance.** Top: Simulated performance (Wikitext-103 validation perplexity (PPL)) versus percent mean-relative simulated noise in feed-forward (FF) and attention (Attn) layers. Systematic errors from experimental data marked with a star. Bottom: comparison of simulated noise model to error from experimental data. The Gaussian shape of the simulated error behavior matches experiment accurately.

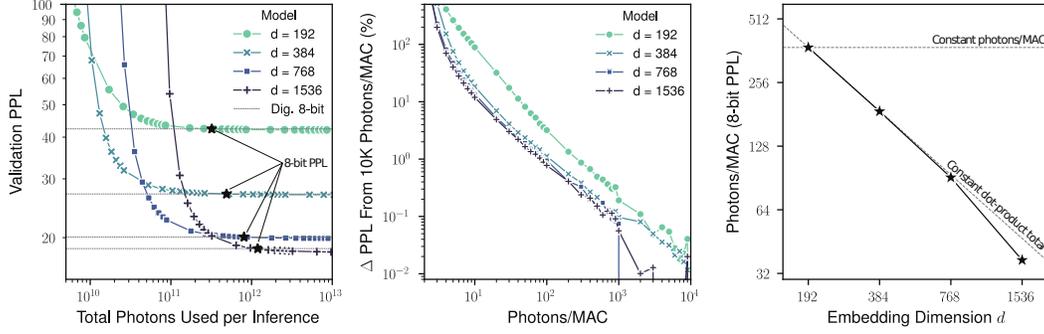


Figure 4: **Simulations of Optical Transformer behavior with varying photon usage.** Left: Wikitext-103 validation-set perplexity (PPL) versus embedding dimension  $d$  and total photons used for a single forward pass/inference. 8-bit digital model performance is shown with dashed lines. Middle: perplexity degrades from ideal with fewer photons-per-MAC; the plot exhibits truncated power-law scaling. Right: Scaling of number of photons needed for an Optical Transformer to achieve the same perplexity as an 8-bit digital-electronic processor, versus model size.

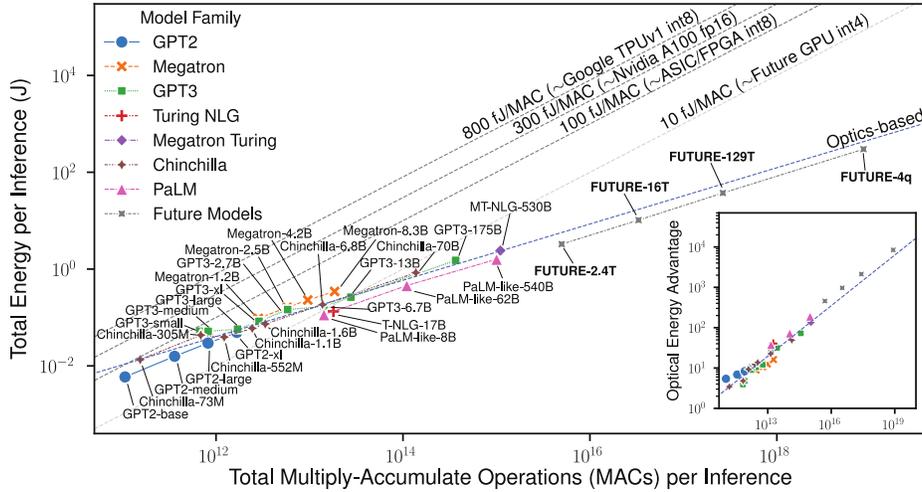


Figure 5: **Estimated energy usage of Transformer models on optical hardware for a single forward pass/inference.** Hypothetical future model designs are labelled **FUTURE-\***. Estimated energy/MAC for digital systems is based on [47]. Trend for energy usage in optical systems (blue) computed based on real models only. Inset: energy advantage of running on optics over estimated NVIDIA A100 usage. The advantage grows with the model compute.  $M = 10^6$ ,  $G = 10^9$ ,  $T = 10^{12}$ ,  $q = 10^{15}$  parameters.

### 285 4.3 Estimated Energy Usage

286 The efficient photon scaling trend we observed in Section 4.2 suggests that Transformers running  
 287 on optical hardware could achieve significant energy efficiency advantages over running on digital  
 288 hardware. To understand the efficiency of Transformers on optical hardware, we designed an ONN  
 289 system based on current hardware that is like our experimental setup, with our measured precision  
 290 and photon scaling. It is an inference system with in-place weights which are loaded once and reused  
 291 forever, activations read from and written to SRAM for every layer, a 10 GHz light modulator array,  
 292 and an optical “core” which can perform 10M multiplications per cycle (this can be thought of as a  
 293 10 megapixel SLM). The photon-per-MAC scaling versus model dimension is taken to be the  $1/d$   
 294 scaling which we found was possible in our simulations, and we assumed that the model operates  
 295 with 5-bit input precision, 8-bit weight precision, and 7-bit output precision, as determined by our  
 296 study of low precision performance in Appendix C. We then calculated according to the approach  
 297 in Section 2.3. For electrical energy we assumed in-place weights and did not include the energy  
 298 for loading them. In Appendix D we explain all assumed energy quantities based on contemporary  
 299 hardware.

300 As models grow, running Transformers on optical hardware has a large and asymptotic efficiency  
301 advantage over running on digital hardware. In Figure 5 we chart estimates of the forward pass energy  
302 required for various models<sup>1</sup>, including a hypothetical family of large, dense Transformer models  
303 designed in a similar fashion, which we label **FUTURE-\***. For comparison, we also chart various  
304 digital systems [47] in different performance regimes, and a hypothetical “next generation” GPU  
305 that can use  $\sim 10$  fJ/MAC. For small models, the optics-based system uses about the same energy,  
306 but eventually gains an advantage that scales asymptotically with the number of MACs. For the  
307 larger models, MT-NLG-530B and FUTURE-4q, the optics-based approach would have  $\sim 140\times$  and  
308  $\sim 8500\times$  energy advantages over the current state-of-the-art GPU (NVIDIA A100) respectively.

309 The breakdown of compute and energy costs by source is in Appendix E. In summary we found that  
310 as models get larger the feed-forward layers require most of the computation, but that the energy of  
311 data access in attention is still very expensive due to the many heads. This is because of the parallel  
312 operation of the Transformer, where the linear layer weights can be re-used for many tokens at a time  
313 (weights-in-place is not possible for attention, and there are  $h n \times n$  attention maps to store).<sup>2</sup>

## 314 5 Discussion

315 The results given in Section 4.3 on optical Transformers’ efficiency have implications for the design  
316 of future ONN hardware/software systems.

317 In Appendix G we discuss in detail the specifications for an ONN system to run large Transformers, as  
318 a target for future work in their design. In summary, we found: once matrix-matrix product operands  
319 exceed  $10^4 \times 10^4$  in size the advantage is significant, and therefore a future ONN should implement at  
320 least this level of parallelism to achieve  $>100\times$  efficiency improvements over current state-of-the-art  
321 GPUs (NVIDIA A100). Given the assumptions we made about weight-maintenance costs in making  
322 our estimates (5.6  $\mu$ W per weight; see Appendix D), an Optical Transformer would need to operate in  
323 the regime where a single matrix-vector multiplication is performed every 0.1 nanoseconds. Current  
324 ONN prototypes either operate at low clock rate or at small scale. Thus building a full ONN system  
325 that realizes the potential benefit is still an open challenge.

326 Future improvements in CMOS technology will be greatly beneficial. In Appendix F we estimate  
327 that future optics-based systems might achieve energy advantages of  $>100,000\times$  running models  
328 the size of FUTURE-4q (over 300 fJ/MAC).

329 Our studies on Transformers illustrates more broadly the relationships between model design and  
330 ONN efficiency. Transformers sought to make large models run efficiently by exploiting hardware’s  
331 strengths in performing large, parallel, dense calculations, and improved in this aspect as they scaled.  
332 As a consequence, as Transformers continue to be optimized for parallel digital electronic hardware,  
333 they will continue to become even more efficient on optical hardware. More generally, architectures  
334 that perform more computations per data access (such as those focusing strongly on linear operations  
335 [58, 35]) will be most promising for optical implementation.

336 **Conclusion** We have demonstrated the ability of Transformer models to run accurately and effi-  
337 ciently on optical hardware through optical experiments and an experiment-informed simulation of  
338 the hardware. We examined Transformers’ scaling behavior with optics and used our findings to  
339 show that optical systems could have a large and asymptotic energy advantage over digital ones that  
340 grows with the model size. For example, we showed that optical hardware may achieve an over  $100\times$   
341 energy advantage when running the largest Transformer models today ( $\sim 500$  billion parameters) and  
342 that larger, future Transformers ( $\sim 4$  quadrillion parameters) may be realized with an  $>8000\times$  optical  
343 energy advantage. We believe our findings about the potential energy-efficiency of optical accelerator  
344 hardware strongly motivate the development of optical processors for large-scale deep learning with  
345 Transformers.

---

<sup>1</sup>The recent PaLM [9] models used a modified architecture. For simpler comparison, we make our estimates using a model with GPT-like architecture but with the PaLM model dimensions, which we call PaLM-Like.

<sup>2</sup>Trends in the design of real models have increasingly favored optics over time. Specifically, attention loads/stores a  $n \times n$  attention matrix for each of the  $h$  attention heads. Models with more MLP compute per attention head have a larger overall ratio of computation to energy usage; larger  $\frac{d}{h}$  is more efficient. The largest GPT2 [43] uses  $\frac{d}{h} = 64$ ; GPT3 [6], 128; MT-NLG-530b [54], 160; and PaLM [9], 384.

## References

- 346
- 347 [1] Michael Andersch, Greg Palmer, Ronny Krashinsky, Nick Stam, Vishal Mehta,  
348 Gonzalo Brito, and Sridhar Ramaswamy. NVIDIA Hopper architecture in-depth.  
349 Technical report, March 2022. URL [https://developer.nvidia.com/blog/  
350 nvidia-hopper-architecture-in-depth/](https://developer.nvidia.com/blog/nvidia-hopper-architecture-in-depth/).
- 351 [2] William Andregg, Michael Andregg, Robert T Weverka, and Lionel Clermont. Wavelength  
352 multiplexed matrix-matrix multiplier, April 19 2019. URL [https://patents.google.com/  
353 patent/US10274989B2/en](https://patents.google.com/patent/US10274989B2/en). (U.S. Patent No. 10,274,989). U.S. Patent and Trademark Office.
- 354 [3] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E. Hinton. Layer normalization, 2016. URL  
355 <https://arxiv.org/abs/1607.06450>.
- 356 [4] Wim Bogaerts, Daniel Pérez, José Capmany, David A B Miller, Joyce Poon, Dirk Englund,  
357 Francesco Morichetti, and Andrea Melloni. Programmable photonic circuits. *Nature*, 586  
358 (7828):207–216, 2020. URL <https://doi.org/10.1038/s41586-020-2764-0>.
- 359 [5] Yelysei Bondarenko, Markus Nagel, and Tijmen Blankevoort. Understanding and overcoming  
360 the challenges of efficient Transformer quantization. In *Proceedings of the 2021 Conference  
361 on Empirical Methods in Natural Language Processing*, pages 7947–7969, Online and Punta  
362 Cana, Dominican Republic, November 2021. Association for Computational Linguistics. URL  
363 <https://aclanthology.org/2021.emnlp-main.627>.
- 364 [6] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal,  
365 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel  
366 Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M.  
367 Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz  
368 Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec  
369 Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners, 2020.  
370 URL <https://arxiv.org/abs/2005.14165>.
- 371 [7] H John Caulfield and Shlomi Dolev. Why future supercomputing requires optics. *Nature  
372 Photonics*, 4(5):261–263, 2010. URL <https://doi.org/10.1038/nphoton.2010.94>.
- 373 [8] Cerebras Systems. Cerebras systems: Achieving industry best AI perfor-  
374 mance through a systems approach. Technical report, Apr 2021. URL [https://  
375 //8968533.fs1.hubspotusercontent-na1.net/hubfs/8968533/whitepapers/  
376 Cerebras-CS-2-Whitepaper.pdf](https://8968533.fs1.hubspotusercontent-na1.net/hubfs/8968533/whitepapers/Cerebras-CS-2-Whitepaper.pdf).
- 377 [9] Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam  
378 Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker  
379 Schuh, Kensen Shi, Sasha Tsvyashchenko, Joshua Maynez, Abhishek Rao, Parker Barnes,  
380 Yi Tay, Noam Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan Du, Ben Hutchinson,  
381 Reiner Pope, James Bradbury, Jacob Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin,  
382 Toju Duke, Anselm Levskaya, Sanjay Ghemawat, Sunipa Dev, Henryk Michalewski, Xavier  
383 Garcia, Vedant Misra, Kevin Robinson, Liam Fedus, Denny Zhou, Daphne Ippolito, David  
384 Luan, Hyeontaek Lim, Barret Zoph, Alexander Spiridonov, Ryan Sepassi, David Dohan, Shivani  
385 Agrawal, Mark Omernick, Andrew M. Dai, Thanumalayan Sankaranarayanan Pillai, Marie Pellat,  
386 Aitor Lewkowycz, Erica Moreira, Rewon Child, Oleksandr Polozov, Katherine Lee, Zongwei  
387 Zhou, Xuezhi Wang, Brennan Saeta, Mark Diaz, Orhan Firat, Michele Catasta, Jason Wei,  
388 Kathy Meier-Hellstern, Douglas Eck, Jeff Dean, Slav Petrov, and Noah Fiedel. PaLM: Scaling  
389 language modeling with Pathways, 2022. URL <https://arxiv.org/abs/2204.02311>.
- 390 [10] Aidan Clark, Diego De Las Casas, Aurelia Guy, Arthur Mensch, Michela Paganini, Jordan  
391 Hoffmann, Bogdan Damoc, Blake Hechtman, Trevor Cai, Sebastian Borgeaud, George Bm  
392 Van Den Driessche, Eliza Rutherford, Tom Hennigan, Matthew J Johnson, Albin Cassirer,  
393 Chris Jones, Elena Buchatskaya, David Budden, Laurent Sifre, Simon Osindero, Oriol  
394 Vinyals, Marc Aurelio Ranzato, Jack Rae, Erich Elsen, Koray Kavukcuoglu, and Karen  
395 Simonyan. Unified scaling laws for routed language models. In Kamalika Chaudhuri,

- 396 Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu, and Sivan Sabato, editors, *Pro-*  
397 *ceedings of the 39th International Conference on Machine Learning*, volume 162 of *Pro-*  
398 *ceedings of Machine Learning Research*, pages 4057–4086. PMLR, 17–23 Jul 2022. URL  
399 <https://proceedings.mlr.press/v162/clark22a.html>.
- 400 [11] Tim Dettmers and Luke Zettlemoyer. The case for 4-bit precision: k-bit inference scaling laws,  
401 2022.
- 402 [12] Tim Dettmers, Mike Lewis, Younes Belkada, and Luke Zettlemoyer. LLM.int8(): 8-bit matrix  
403 multiplication for Transformers at scale, 2022. URL <https://arxiv.org/abs/2208.07339>.
- 404 [13] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of  
405 deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Confer-*  
406 *ence of the North American Chapter of the Association for Computational Linguistics: Human*  
407 *Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis,  
408 Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1423.  
409 URL <https://aclanthology.org/N19-1423>.
- 410 [14] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai,  
411 Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly,  
412 Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image  
413 recognition at scale. *ICLR*, 2021.
- 414 [15] Johannes Feldmann, Nathan Youngblood, Maxim Karpov, Helge Gehring, Xuan Li, Maik  
415 Stappers, Manuel Le Gallo, Xin Fu, Anton Lukashchuk, Arslan Sajid Raja, et al. Parallel  
416 convolutional processing using an integrated photonic tensor core. *Nature*, 589(7840):52–58,  
417 2021.
- 418 [16] Graphcore. The data center architecture for graphcore computing. Techni-  
419 cal report, Apr 2021. URL [https://www.graphcore.ai/hubfs/](https://www.graphcore.ai/hubfs/Graphcore-Mk2-IPU-System-Architecture-GC.pdf)  
420 [Graphcore-Mk2-IPU-System-Architecture-GC.pdf](https://www.graphcore.ai/hubfs/Graphcore-Mk2-IPU-System-Architecture-GC.pdf).
- 421 [17] Habana Labs. HABANA® GAUDI®2 white paper. Technical report, June 2022. URL  
422 <https://habana.ai/wp-content/uploads/pdf/2022/audi2-whitepaper.pdf>.
- 423 [18] Ryan Hamerly, Liane Bernstein, Alexander Sludds, Marin Soljačić, and Dirk Englund. Large-  
424 scale optical neural networks based on photoelectric multiplication. *Physical Review X*, 9(2):  
425 021032, 2019. URL <https://doi.org/10.1103/PhysRevX.9.021032>.
- 426 [19] Yoshio Hayasaki, Ichiro Tohyama, Toyohiko Yatagai, Masahiko Mori, and Satoshi Ishihara.  
427 Optical learning neural network using Selfoc microlens array. *Japanese Journal of Applied*  
428 *Physics*, 31(5S):1689, 1992. URL <https://doi.org/10.1143/JJAP.31.1689>.
- 429 [20] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image  
430 recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.  
431 IEEE, June 2016. doi: 10.1109/cvpr.2016.90. URL [https://doi.org/10.1109/cvpr.2016.](https://doi.org/10.1109/cvpr.2016.90)  
432 90.
- 433 [21] Dan Hendrycks and Kevin Gimpel. Gaussian error linear units (GELUs), 2016. URL <https://arxiv.org/abs/1606.08415>.
- 435 [22] Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza  
436 Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, Tom  
437 Hennigan, Eric Noland, Katie Millican, George van den Driessche, Bogdan Damoc, Aurelia  
438 Guy, Simon Osindero, Karen Simonyan, Erich Elsen, Jack W. Rae, Oriol Vinyals, and Laurent  
439 Sifre. Training compute-optimal large language models, 2022. URL [https://arxiv.org/](https://arxiv.org/abs/2203.15556)  
440 [abs/2203.15556](https://arxiv.org/abs/2203.15556).
- 441 [23] Andrew G. Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias  
442 Weyand, Marco Andreetto, and Hartwig Adam. MobileNets: Efficient convolutional neural  
443 networks for mobile vision applications. *ArXiv*, abs/1704.04861, 2017.

- 444 [24] Chaoran Huang, Volker J. Sorger, Mario Miscuglio, Mohammed Al-Qadasi, Avilash Mukherjee,  
445 Lutz Lampe, Mitchell Nichols, Alexander N. Tait, Thomas Ferreira de Lima, Bicky A. Marquez,  
446 Jiahui Wang, Lukas Chrostowski, Mable P. Fok, Daniel Brunner, Shanhui Fan, Sudip Shekhar,  
447 Paul R. Prucnal, and Bhavin J. Shastri. Prospects and applications of photonic neural networks.  
448 *Advances in Physics: X*, 7(1), October 2021. doi: 10.1080/23746149.2021.1981155. URL  
449 <https://doi.org/10.1080/23746149.2021.1981155>.
- 450 [25] Benoit Jacob, Skirmantas Kligys, Bo Chen, Menglong Zhu, Matthew Tang, Andrew Howard,  
451 Hartwig Adam, and Dmitry Kalenichenko. Quantization and training of neural networks for  
452 efficient integer-arithmetic-only inference. In *Proceedings of the IEEE Conference on Computer  
453 Vision and Pattern Recognition*, pages 2704–2713, 2018. URL [https://doi.org/10.1109/  
454 CVPR.2018.00286](https://doi.org/10.1109/CVPR.2018.00286).
- 455 [26] Andrew Jaegle, Sebastian Borgeaud, Jean-Baptiste Alayrac, Carl Doersch, Catalin Ionescu,  
456 David Ding, Skanda Koppula, Daniel Zoran, Andrew Brock, Evan Shelhamer, Olivier Hénaff,  
457 Matthew M. Botvinick, Andrew Zisserman, Oriol Vinyals, and João Carreira. Perceiver IO: A  
458 general architecture for structured inputs & outputs, 2021. URL [https://arxiv.org/abs/  
459 2107.14795](https://arxiv.org/abs/2107.14795).
- 460 [27] Andrew Jaegle, Felix Gimeno, Andy Brock, Oriol Vinyals, Andrew Zisserman, and Joao  
461 Carreira. Perceiver: General perception with iterative attention. In Marina Meila and Tong  
462 Zhang, editors, *Proceedings of the 38th International Conference on Machine Learning*, volume  
463 139 of *Proceedings of Machine Learning Research*, pages 4651–4664. PMLR, 18–24 Jul 2021.  
464 URL <https://proceedings.mlr.press/v139/jaegle21a.html>.
- 465 [28] Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B. Brown, Benjamin Chess, Rewon Child,  
466 Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language  
467 models, 2020. URL <https://arxiv.org/abs/2001.08361>.
- 468 [29] Hyungjun Kim, Jihoon Park, Changhun Lee, and Jae-Joon Kim. Improving accuracy of binary  
469 neural networks using unbalanced activation distribution. In *2021 IEEE/CVF Conference on  
470 Computer Vision and Pattern Recognition (CVPR)*, pages 7858–7867, 2021. doi: 10.1109/  
471 CVPR46437.2021.00777.
- 472 [30] Jinwoo Kim, Tien Dat Nguyen, Seonwoo Min, Sungjun Cho, Moontae Lee, Honglak Lee, and  
473 Seunghoon Hong. Pure Transformers are powerful graph learners. *arXiv*, abs/2207.02505, 2022.  
474 URL <https://arxiv.org/abs/2207.02505>.
- 475 [31] Alex Krizhevsky. Convolutional deep belief networks on cifar-10. 2010. URL [https:  
476 //www.cs.toronto.edu/~kriz/conv-cifar10-aug2010.pdf](https://www.cs.toronto.edu/~kriz/conv-cifar10-aug2010.pdf).
- 477 [32] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with  
478 deep convolutional neural networks. In F. Pereira, C.J. Burges, L. Bottou, and K.Q. Wein-  
479 berger, editors, *Advances in Neural Information Processing Systems*, volume 25. Curran  
480 Associates, Inc., 2012. URL [https://proceedings.neurips.cc/paper/2012/file/  
481 c399862d3b9d6b76c8436e924a68c45b-Paper.pdf](https://proceedings.neurips.cc/paper/2012/file/c399862d3b9d6b76c8436e924a68c45b-Paper.pdf).
- 482 [33] Aitor Lewkowycz, Anders Andreassen, David Dohan, Ethan Dyer, Henryk Michalewski, Vinay  
483 Ramasesh, Ambrose Slone, Cem Anil, Imanol Schlag, Theo Gutman-Solo, Yuhuai Wu, Behnam  
484 Neyshabur, Guy Gur-Ari, and Vedant Misra. Solving quantitative reasoning problems with  
485 language models, 2022. URL <https://arxiv.org/abs/2206.14858>.
- 486 [34] Zhuohan Li, Eric Wallace, Sheng Shen, Kevin Lin, Kurt Keutzer, Dan Klein, and Joseph E.  
487 Gonzalez. Train large, then compress: Rethinking model size for efficient training and inference  
488 of transformers. In *Proceedings of the 37th International Conference on Machine Learning*,  
489 ICML’20. JMLR.org, 2020.
- 490 [35] Hanxiao Liu, Zihang Dai, David R. So, and Quoc V. Le. Pay attention to mlps, 2021. URL  
491 <https://arxiv.org/abs/2105.08050>.
- 492 [36] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining  
493 Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings  
494 of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021.

- 495 [37] Kevin Lu, Aditya Grover, Pieter Abbeel, and Igor Mordatch. Pretrained transformers as universal  
496 computation engines. *arXiv preprint arXiv:2103.05247*, 2021.
- 497 [38] Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. Pointer sentinel mixture  
498 models. In *International Conference on Learning Representations (ICLR)*, 2017. URL <https://openreview.net/forum?id=Byj72udxe>.  
499
- 500 [39] Charis Mesaritakis, Vassilis Papataxiarhis, and Dimitris Syvridis. Micro ring resonators as  
501 building blocks for an all-optical high-speed reservoir-computing bit-pattern-recognition system.  
502 *J. Opt. Soc. Am. B*, 30(11):3048–3055, Nov 2013. doi: 10.1364/JOSAB.30.003048. URL  
503 <https://opg.optica.org/josab/abstract.cfm?URI=josab-30-11-3048>.
- 504 [40] Mario Miscuglio, Zibo Hu, Shurui Li, Jonathan K George, Roberto Capanna, Hamed Dalir,  
505 Philippe M Bardet, Puneet Gupta, and Volker J. Sorger. Massively parallel amplitude-only  
506 fourier neural network. *Optica*, 7(12):1812–1819, 2020. URL <https://doi.org/10.1364/OPTICA.408659>.  
507
- 508 [41] Mitchell A Nahmias, Thomas Ferreira De Lima, Alexander N Tait, Hsuan-Tung Peng, Bhavin J  
509 Shastri, and Paul R Prucnal. Photonic multiply-accumulate operations for neural networks.  
510 *IEEE Journal of Selected Topics in Quantum Electronics*, 26:1–18, 2020. URL <https://doi.org/10.1109/JSTQE.2019.2941485>.  
511
- 512 [42] Alec Radford and Karthik Narasimhan. Improving language understanding by generative  
513 pre-training. 2018. URL <https://openai.com/blog/language-unsupervised/>.
- 514 [43] Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Lan-  
515 guage models are unsupervised multitask learners. 2019. URL <https://openai.com/blog/better-language-models/>.  
516
- 517 [44] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agar-  
518 wal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya  
519 Sutskever. Learning transferable visual models from natural language supervision. In Marina  
520 Meila and Tong Zhang, editors, *Proceedings of the 38th International Conference on Machine  
521 Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 8748–8763. PMLR,  
522 18–24 Jul 2021. URL <https://proceedings.mlr.press/v139/radford21a.html>.
- 523 [45] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark  
524 Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In Marina Meila and Tong Zhang,  
525 editors, *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of  
526 *Proceedings of Machine Learning Research*, pages 8821–8831. PMLR, 18–24 Jul 2021. URL  
527 <https://proceedings.mlr.press/v139/ramesh21a.html>.
- 528 [46] Scott Reed, Konrad Zolna, Emilio Parisotto, Sergio Gómez Colmenarejo, Alexander Novikov,  
529 Gabriel Barth-maroon, Mai Giménez, Yury Sulsky, Jackie Kay, Jost Tobias Springenberg, Tom  
530 Eccles, Jake Bruce, Ali Razavi, Ashley Edwards, Nicolas Heess, Yutian Chen, Raia Hadsell,  
531 Oriol Vinyals, Mahyar Bordbar, and Nando de Freitas. A generalist agent. *Transactions on  
532 Machine Learning Research*, 2022. ISSN 2835-8856. URL [https://openreview.net/  
533 forum?id=1ikK0kHjvj](https://openreview.net/forum?id=1ikK0kHjvj). Featured Certification.
- 534 [47] Albert Reuther, Peter Michaleas, Michael Jones, Vijay Gadepally, Siddharth Samsi, and Jeremy  
535 Kepner. Survey of machine learning accelerators. *arXiv:2009.00993*, 2020. URL <https://arxiv.org/abs/2009.00993>.  
536
- 537 [48] Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. DistilBERT, a distilled  
538 version of BERT: smaller, faster, cheaper and lighter. *ArXiv*, abs/1910.01108, 2019.
- 539 [49] Abu Sebastian, Manuel Le Gallo, Riduan Khaddam-Aljameh, and Evangelos Eleftheriou.  
540 Memory devices and applications for in-memory computing. *Nature Nanotechnology*, 15(7):  
541 529–544, March 2020. doi: 10.1038/s41565-020-0655-z. URL [https://doi.org/10.1038/  
542 s41565-020-0655-z](https://doi.org/10.1038/s41565-020-0655-z).

- 543 [50] Jaime Sevilla, Lennart Heim, Anson Ho, Tamay Besiroglu, Marius Hobbhahn, and Pablo  
544 Villalobos. Compute trends across three eras of machine learning. In *2022 International Joint  
545 Conference on Neural Networks (IJCNN)*, pages 1–8, 2022. doi: 10.1109/IJCNN55064.2022.  
546 9891914.
- 547 [51] Bhavin J Shastri, Alexander N Tait, T Ferreira de Lima, Wolfram HP Pernice, Harish Bhaskaran,  
548 C David Wright, and Paul R Prucnal. Photonics for artificial intelligence and neuromorphic  
549 computing. *Nature Photonics*, 15(2):102–114, 2021. URL [https://doi.org/10.1038/  
550 s41566-020-00754-y](https://doi.org/10.1038/s41566-020-00754-y).
- 551 [52] Yichen Shen, Nicholas C Harris, Scott Skirlo, Mihika Prabhu, Tom Baehr-Jones, Michael  
552 Hochberg, Xin Sun, Shijie Zhao, Hugo Larochelle, Dirk Englund, and Marin Soljačić. Deep  
553 learning with coherent nanophotonic circuits. *Nature Photonics*, 11(7):441, 2017. URL  
554 <https://doi.org/10.1038/nphoton.2017.93>.
- 555 [53] Alexander Sludds, Saumil Bandyopadhyay, Zaijun Chen, Zhizhen Zhong, Jared Cochrane,  
556 Liane Bernstein, Darius Bunandar, P. Ben Dixon, Scott A. Hamilton, Matthew Streshinsky,  
557 Ari Novack, Tom Baehr-Jones, Michael Hochberg, Manya Ghobadi, Ryan Hamerly, and Dirk  
558 Englund. Delocalized photonic deep learning on the internet’s edge. *Science*, 378(6617):270–  
559 276, 2022. doi: 10.1126/science.abq8271. URL [https://www.science.org/doi/abs/10.  
560 1126/science.abq8271](https://www.science.org/doi/abs/10.1126/science.abq8271).
- 561 [54] Shaden Smith, Mostofa Patwary, Brandon Norick, Patrick LeGresley, Samyam Rajbhandari,  
562 Jared Casper, Zhun Liu, Shrimai Prabhumoye, George Zerveas, Vijay Korthikanti, Elton Zhang,  
563 Rewon Child, Reza Yazdani Aminabadi, Julie Bernauer, Xia Song, Mohammad Shoeybi,  
564 Yuxiong He, Michael Houston, Saurabh Tiwary, and Bryan Catanzaro. Using DeepSpeed and  
565 Megatron to train Megatron-Turing NLG 530B, a large-scale generative language model, 2022.  
566 URL <https://arxiv.org/abs/2201.11990>.
- 567 [55] James Spall, Xianxin Guo, Thomas D Barrett, and AI Lvovsky. Fully reconfigurable coherent  
568 optical vector–matrix multiplication. *Optics Letters*, 45(20):5752–5755, 2020. URL <https://doi.org/10.1364/OL.401675>.  
569
- 570 [56] Pascal Stark, Folkert Horst, Roger Dangel, Jonas Weiss, and Bert Jan Offrein. Opportunities for  
571 integrated photonic neural networks. *Nanophotonics*, 9(13):4221–4232, 2020. URL <https://doi.org/10.1515/nanoph-2020-0297>.  
572
- 573 [57] Alexander N. Tait, John Chang, Bhavin J. Shastri, Mitchell A. Nahmias, and Paul R. Prucnal.  
574 Demonstration of WDM weighted addition for principal component analysis. *Opt. Express*,  
575 23(10):12758–12765, May 2015. doi: 10.1364/OE.23.012758. URL [https://opg.optica.  
576 org/oe/abstract.cfm?URI=oe-23-10-12758](https://opg.optica.org/oe/abstract.cfm?URI=oe-23-10-12758).
- 577 [58] Ilya Tolstikhin, Neil Houlsby, Alexander Kolesnikov, Lucas Beyer, Xiaohua Zhai, Thomas  
578 Unterthiner, Jessica Yung, Andreas Steiner, Daniel Keysers, Jakob Uszkoreit, Mario Lucic,  
579 and Alexey Dosovitskiy. MLP-Mixer: An all-MLP architecture for vision. *arXiv preprint  
580 arXiv:2105.01601*, 2021.
- 581 [59] Marcos Treviso, Tianchu Ji, Ji-Ung Lee, Betty van Aken, Qingqing Cao, Manuel R. Ciosici,  
582 Michael Hassid, Kenneth Heafield, Sara Hooker, Pedro H. Martins, André F. T. Martins,  
583 Peter Milder, Colin Raffel, Edwin Simpson, Noam Slonim, Niranjan Balasubramanian, Leon  
584 Derczynski, and Roy Schwartz. Efficient methods for natural language processing: A survey,  
585 2022. URL <https://arxiv.org/abs/2209.00099>.
- 586 [60] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N  
587 Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In  
588 *Advances in Neural Information Processing Systems*, pages 5998–6008. Curran As-  
589 sociates, Inc., 2017. URL [https://proceedings.neurips.cc/paper/2017/file/  
590 3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf](https://proceedings.neurips.cc/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf).
- 591 [61] Tianyu Wang, Shi-Yuan Ma, Logan G. Wright, Tatsuhiko Onodera, Brian C. Richard, and  
592 Peter L. McMahon. An optical neural network using less than 1 photon per multiplication.  
593 *Nature Communications*, 13(1), January 2022. doi: 10.1038/s41467-021-27774-8. URL  
594 <https://doi.org/10.1038/s41467-021-27774-8>.

- 595 [62] Gordon Wetzstein, Aydogan Ozcan, Sylvain Gigan, Shanhui Fan, Dirk Englund, Marin Soljačić,  
596 Cornelia Denz, David A. B. Miller, and Demetri Psaltis. Inference in artificial intelligence  
597 with deep optics and photonics. *Nature*, 588(7836):39–47, Dec 2020. ISSN 1476-4687. doi:  
598 10.1038/s41586-020-2973-6. URL <https://doi.org/10.1038/s41586-020-2973-6>.
- 599 [63] Changming Wu, Heshan Yu, Seokhyeong Lee, Ruoming Peng, Ichiro Takeuchi, and Mo Li.  
600 Programmable phase-change metasurfaces on waveguides for multimode photonic convolutional  
601 neural network. *arXiv preprint arXiv:2004.10651*, 2020.
- 602 [64] Xingyuan Xu, Mengxi Tan, Bill Corcoran, Jiayang Wu, Andreas Boes, Thach G Nguyen, Sai T  
603 Chu, Brent E Little, Damien G Hicks, Roberto Morandotti, Arnan Mitchell, and David J Moss.  
604 11 TOPS photonic convolutional accelerator for optical neural networks. *Nature*, 589(7840):  
605 44–51, 2021. URL <https://doi.org/10.1038/s41586-020-03063-0>.
- 606 [65] Jiahui Yu, Zirui Wang, Vijay Vasudevan, Legg Yeung, Mojtaba Seyedhosseini, and Yonghui  
607 Wu. CoCa: Contrastive captioners are image-text foundation models. *Transactions on Machine  
608 Learning Research*, 2022. ISSN 2835-8856. URL [https://openreview.net/forum?id=  
609 Ee277P3AYC](https://openreview.net/forum?id=Ee277P3AYC).
- 610 [66] Xiaohua Zhai, Alexander Kolesnikov, Neil Houlsby, and Lucas Beyer. Scaling vision Transform-  
611 ers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition  
612 (CVPR)*, pages 12104–12113, June 2022.