

A Mathematical Explanation of Transformers for Large Language Models and GPTs *

Xue-Cheng Tai[†], Hao Liu[‡], Lingfeng Li[§], and Raymond H. Chan[¶]

Abstract. The Transformer architecture has revolutionized the field of sequence modeling and underpins the recent breakthroughs in large language models (LLMs). However, a comprehensive mathematical theory that explains its structure and operations remains elusive. In this work, we propose a novel continuous framework that rigorously interprets the Transformer as a discretization of a structured integro-differential equation. Within this formulation, the self-attention mechanism emerges naturally as a non-local integral operator, and layer normalization is characterized as a projection to a time-dependent constraint. This operator-theoretic and variational perspective offers a unified and interpretable foundation for understanding the architecture’s core components, including attention, feedforward layers, and normalization. Our approach extends beyond previous theoretical analyses by embedding the entire Transformer operation in continuous domains for both token indices and feature dimensions. This leads to a principled and flexible framework that not only deepens theoretical insight but also offers new directions for architecture design, analysis, and control-based interpretations. This new interpretation provides a step toward bridging the gap between deep learning architectures and continuous mathematical modeling, and contributes a foundational perspective to the ongoing development of interpretable and theoretically grounded neural network models.

Key words. Transformer, attention, operator splitting, integro-differential equation

MSC codes. 68U10, 94A08

1. Introduction. Deep neural networks (DNNs) have revolutionized numerous fields, including natural language processing [15, 55], computer vision [23], scientific computing [31, 36, 26], and medical diagnostics [35, 22]. Among these architectures, Transformers [46] have recently emerged as particularly powerful tools, underpinning remarkable successes in large language models (LLMs) such as GPT-3 and GPT-4 [46]. Besides language processing [59], the Transformer and its variants have also been applied in applications, including image processing [8, 41], graph processing [56], operator learning [6, 3, 28, 54].

Recently, a series of works have been conducted to study the theoretical foundation or interpretability of Transformers. Approximation and generalization error of Transformers were studied in [10, 44, 19]. It was shown in [19, 40] that Transformers are adaptive to data’s low-dimensional structures. For interpretability, in [24], it was shown that Transformers are cubic or high-order splines. It was shown in [21] that the vision Transformer can learn spatial structures of the dataset. The Transformer was interpreted as an ODE solver for multi-particle dynamical systems in [9, 11, 33]. Meng et al. [?] showed that Transformer is a step for solving

*Submitted to the editors DATE.

Funding:

[†]Norwegian Research Centre, Bergen, Norway (xtai@norceresearch.no, xuechengtai@gmail.com).

[‡]Department of Mathematics, Hong Kong Baptist University, Kowloon Tong, Hong Kong (haoliu@hkbu.edu.hk).

[§]Hong Kong Centre for Cerebro-Cardiovascular Health Engineering, Shatin, Hong Kong (lfl@hkcoche.org)

[¶]School of Data Science and Department of Operations and Risk Management, Lingnan University, Tuen Mun, Hong Kong (raymond.chan@ln.edu.hk).

36 a nonlocal variational model. A general introduction to Transformers with mathematical
37 descriptions can be found in [45].

38 Compared to Transformers, more studies can be found for deep neural networks (DNNs).
39 A widely accepted interpretation views DNN as nonlinear high-dimensional function approx-
40 imators [53, 30, 58, 4]. Under this view, the layers of a DNN perform successive nonlinear
41 transformations that approximate complex mappings from input to output data. This ap-
42 proximation framework has been instrumental in explaining the expressive power of DNNs.
43 Another promising direction for explaining and designing neural network architectures in-
44 volves interpreting them as discretizations of continuous-time dynamical systems governed by
45 differential or integral equations [5, 49, 2, 48]. E et al. show that DNNs can be viewed as
46 continuous dynamical systems in [48] and special discretization of continuous problems in [49].
47 In [5], Chen et al. pointed out connections between residual networks and ordinary differen-
48 tial equations, and proposed Neural ODE. Benning et al. investigated the relations between
49 DNNs and optimal control problems [2]. In this class of approach, neural network training
50 is seen as an optimal control problem constrained by a continuous evolution equation. Re-
51 lated to the above-mentioned approach, we especially want to mention some recent studies
52 that have proposed a general framework for designing neural network architectures inspired
53 by continuous-time dynamical systems [42, 27]. These methods leverage operator-splitting
54 techniques to construct neural networks, which are so-called continuous neural networks, as
55 discretizations of continuous differential equations. The continuous learning problem is an
56 optimal control problem with this "neural network" as a constraint. In this framework, the
57 evolution of hidden states in a network is governed by differential equations (continuous neural
58 network). The goal of this continuous learning problem is to find optimal controls—represented
59 by learnable parameters—within a discretized dynamical process.

60 The use of operator-splitting allows these complex dynamics to be decomposed into sim-
61 pler substeps, each of which corresponds to a specific type of operation (e.g., convolution,
62 nonlinearity, normalization). These substeps are then unrolled into individual layers of a
63 neural network. This interpretation not only clarifies the structure and purpose of different
64 architectural components but also facilitates the incorporation of prior knowledge about the
65 underlying data dynamics or physics.

66 Building on this theory, we have recently provided a rigorous mathematical explanation
67 of the well-known UNet architecture [43]. We demonstrated that UNet can be viewed as a
68 specific discretization of the following simple differential equation, where each encoder and
69 decoder stage corresponds to a particular substep in the operator-splitting scheme:

$$70 \quad (1.1) \quad \begin{cases} \frac{\partial u(\mathbf{x}, t)}{\partial t} = W(\mathbf{x}, t) * u(\mathbf{x}, t) + d(t) - \ln \frac{u(\mathbf{x}, t)}{1-u(\mathbf{x}, t)} + \partial I_{\Sigma}(u), & (\mathbf{x}, t) \in \Omega \times (0, T], \\ u(\mathbf{x}, 0) = f(\mathbf{x}), & \mathbf{x} \in \Omega, \end{cases}$$

71 where $*$ denotes convolution, I_{Σ} is the indicator function of $\Sigma = \{u : u(\mathbf{x}) \geq 0 \text{ for } \mathbf{x} \in \Omega\}$,
72 f is a function defined in a domain Ω and T is the final time chosen by the user for time
73 variable t . The learnable variables $\theta = \{W(\mathbf{x}, t), d(t)\}$ are regarded as control variables which
74 will be learned in the training process. The mapping $\mathcal{N}_{\theta}(f) \mapsto u(\mathbf{x}, T)$ turns out to be the
75 UNet in the continuous setting. This perspective not only elucidates the structure of the UNet
76 architecture but also provides a foundation for systematic modifications and enhancements.

77 A range of alternative interpretations and mathematically inspired designs of deep neural
 78 networks can be found in the literature, including [39, 2, 38, 17, 20, 25, 29, 16, 52, ?]. For
 79 example, neural networks have been interpreted as optimal control problems in [2, 38]. Gregor
 80 and LeCun [16] were among the first to connect deep learning with sparse coding through
 81 learned iterative shrinkage algorithms. Ruthotto and Haber [39] proposed neural networks
 82 derived from partial differential equations (PDEs), while He and Xu [20] introduced MgNet,
 83 which employs multigrid principles for feature extraction. Hagemann et al. [18] formulated
 84 stochastic normalizing flows as Markov chains to tackle inverse problems. Long et al. [29]
 85 built PDENet based on finite difference schemes to recover PDE models from data. Yang
 86 et al. [52] constructed ADMM-CSNet by unrolling the ADMM optimization algorithm into
 87 a network architecture. Wang et al. [47] developed a diffusion-residual network inspired by
 88 convection-diffusion ODEs for classification tasks. Most recently, Martin et al. [34] introduced
 89 a plug-and-play architecture using flow matching for image restoration. Liu et al. [?] proposed
 90 a novel network regularizer inspired by the reverse process of PDE-based evolution models.
 91 Networks regularized by prior information are proposed in [?, ?, ?]. These pioneering efforts
 92 have provided essential insights and have directly inspired our recent work [42, 43] and the
 93 present study.

94 In this work, we will develop a mathematical theoretical framework to explain the Trans-
 95 former architecture of [46], which has achieved remarkable success in large language models
 96 (LLMs) and other sequence modeling tasks. We show that the Transformer can be interpreted
 97 as a discretization of the following continuous integro-differential equation:

$$(1.2) \quad \left\{ \begin{array}{l} u_t = \underbrace{\langle \gamma(\mathbf{x}, \cdot, t; u), V(\cdot, \mathbf{y}, t; u) \rangle_{\Omega_x}}_{\text{I: attention}} + \underbrace{\partial I_{S_1(\sigma_1(t), \sigma_2(t))}(u)}_{\text{II: layer normalization}} \\ \quad + \underbrace{\sum_{j=1}^J \left(\langle W_j(\cdot, \mathbf{y}, t) u(\mathbf{x}, \cdot, t) \rangle_{\Omega_y} + b_j(\mathbf{x}, t) \right)}_{\text{III: fully connected network}} + \partial I_{S_2}(u) \text{ for } t \in (0, T], \\ u(\mathbf{x}, \mathbf{y}, 0) = f(\mathbf{x}, \mathbf{y}), \end{array} \right.$$

99 for $(\mathbf{x}, \mathbf{y}) \in \Omega_x \times \Omega_y$, where \mathbf{x} denotes the token index, \mathbf{y} the entry of token vectors, Ω_x and
 100 Ω_y are continuous domains, and $T > 0$ is a fixed terminal time. This formulation provides
 101 a rigorous mathematical interpretation of the Transformer and its key components, such as
 102 self-attention and feedforward layers. The definitions of attention, layer normalization, and
 103 the fully connected layer, along with detailed explanations of the domains Ω_x , Ω_y , and the
 104 control variables, are given in Section 2. Equation (1.2) highlights the non-local nature of
 105 attention and reveals the Transformer as a structured operator derived from a discretized
 106 variational principle. We note that several previous works have also established connections
 107 between Transformer-based models and multi-particle dynamical systems [9, 11, 33]. These
 108 studies interpret the original or modified Transformer models as discretizations of interacting
 109 particle systems in a time-continuous setting, treating interactions across spatial variables as
 110 particle interactions. In contrast, our formulation (1.2) is fundamentally different: we show

111 that the Transformer arises as a discretization of a continuous integro-differential equation,
112 offering a more unified operator-theoretic and variational perspective.

113 This perspective offers several advantages. First, it provides a unifying mathematical
114 framework that connects diverse architectures, such as CNNs, UNets, and Transformers, un-
115 der the common lens of differential and integral equations. This unification improves our
116 understanding of the design principles underlying modern deep networks and provides insights
117 for cross-architectural studies.

118 Second, by casting neural network architectures as time-stepping schemes of dynamical
119 systems, our framework allows the systematic exploration of new architectures using well-
120 established tools from numerical analysis. For instance, stability, convergence, and approx-
121 imation properties of the underlying continuous models can inform the selection of network
122 structures and hyperparameters, leading to more robust and interpretable models.

123 Third, this approach creates a principled pathway for embedding domain-specific knowl-
124 edge, such as physical laws, geometric structures, or conservation principles, directly into the
125 design of neural architectures, leading to architectures tailored for specific scientific or engi-
126 neering tasks.

127 In summary, our formulation not only advances the theoretical understanding of deep
128 neural networks but also provides actionable tools for the principled design of next-generation
129 models. It bridges the gap between continuous mathematical modeling and discrete network
130 implementation, establishing a foundation for informed, explainable, and application-aware
131 neural network development.

132 This paper is organized as follows: We introduce the integral equation for Transformer
133 with single-head attention in continuous setting in Section 2. An operator-splitting based
134 algorithm is proposed in Section 3 to solve the integral equation. We show in Section 4 that
135 the discretized algorithm exactly recovers the Transformer encoder in [46], and also discuss
136 how to adapt the proposed algorithm to recover the architecture of ViT [8]. In Section 5, we
137 extend the integral equation and proposed algorithm to Transformer with multi-head attention.
138 Extension to convolutional Transformer and its relation to CvT [50] are discussed in Section
139 6. This paper is concluded in Section 7.

140 In the rest of this paper, we use bold lowercase letters to denote vectors, normal letters to
141 denote scalars and functions, and calligraphic letters to denote operators.

142 **2. The Continuous Model.** In the following, we explain the Transformer in the seminal
143 work in [46] in the continuous setting and show later that proper discretizations of the proposed
144 model will recover the Transformer exactly.

145 **2.1. The continuous setting for Transformers.** In this section, we propose a control
146 problem in a continuous setting. The continuous Transformer is a time-dependent integral
147 equation. Later, we will show that after discretizations by operator splitting for the time
148 variable t and spatial discretizations for the \mathbf{x}, \mathbf{y} variables, the resulting splitting method
149 recovers the Transformer architecture. The Transformer proposed in [46] consists of three
150 components: attention, layer normalization, and feedforward networks. In the derivation of our
151 control problem, we show that each component has a corresponding operation in the integral
152 equation. We first focus on single-head attention. Our framework can be easily extended to
153 the multi-head case, which will be detailed in Section 5.

154 To make the explanation clearer and simpler, our notations and setup will follow [46], but
 155 in a continuous counterpart. Let $u(\mathbf{x}, \mathbf{y}, t)$ be a function of \mathbf{x}, \mathbf{y} and time t . In applications of
 156 LLMs, $\mathbf{x} \in \Omega_x$ refers to the index of tokens, $\mathbf{y} \in \Omega_y$ refers to entries of token vectors. Please
 157 refer to Appendix A for explanations of the physical meanings of \mathbf{x}, \mathbf{y} , and u for applications
 158 of Transformers to LLM.

159 As stated in [46], attention is all you need for the Transformer. In fact, we show below that
 160 the attention layer has a very simple mathematical explanation in the continuous setting. It
 161 consists of integral transformations on the input functions to extract features with three differ-
 162 ent kernels. We use features extracted from the first two kernels to generate the attention score
 163 through a softmax operator, and then multiply the attention score by the feature extracted
 164 from the third kernel to compute the output. All three kernels in the integral transformation
 165 will be parameterized as learnable variables, see (2.11).

166 **2.2. Integral Transformations for Attention and Feature Extraction.** Let the functions
 167 $W^Q(\mathbf{y}, \tilde{\mathbf{y}}, t)$, $W^K(\mathbf{y}, \tilde{\mathbf{y}}, t)$, and $W^V(\mathbf{y}, \tilde{\mathbf{y}}, t)$ be three kernels defined on $\Omega_y \times \Omega_y \times [0, T]$, and
 168 they will be learned from data as shown in LLM applications. See Appendix A for explanations
 169 of the physical meanings of $\mathbf{x}, \mathbf{y}, u(\mathbf{x}, \mathbf{y}, t)$. For a given function $u(\mathbf{x}, \mathbf{y}, t)$, we define integral
 170 transformations:

$$171 \quad (2.1) \quad Q(\mathbf{x}, \mathbf{y}, t; u) = \left\langle W^Q(\cdot, \mathbf{y}, t), u(\mathbf{x}, \cdot, t) \right\rangle_{\Omega_y} = \int_{\Omega_y} W^Q(\xi, \mathbf{y}, t) u(\mathbf{x}, \xi, t) d\xi,$$

$$172 \quad (2.2) \quad K(\mathbf{x}, \mathbf{y}, t; u) = \left\langle W^K(\cdot, \mathbf{y}, t), u(\mathbf{x}, \cdot, t) \right\rangle_{\Omega_y} = \int_{\Omega_y} W^K(\xi, \mathbf{y}, t) u(\mathbf{x}, \xi, t) d\xi,$$

$$173 \quad (2.3) \quad V(\mathbf{x}, \mathbf{y}, t; u) = \left\langle W^V(\cdot, \mathbf{y}, t), u(\mathbf{x}, \cdot, t) \right\rangle_{\Omega_y} = \int_{\Omega_y} W^V(\xi, \mathbf{y}, t) u(\mathbf{x}, \xi, t) d\xi.$$

174 Here and later, we use $\langle \cdot, \cdot \rangle_{\Omega_y}$ to denote L^2 inner product on a given domain Ω_y . Its induced
 175 norm will be denoted by $\| \cdot \|_{\Omega_y}$. Functions Q and K will be used to generate attention scores,
 176 and V is used to extract features from u . When all these functions are discretized over a
 177 pixel-type grid over $\Omega_x \times \Omega_y$, the discrete counterparts of these functions are represented by
 178 matrices and the discrete counterparts of the integral transformations are just standard matrix
 179 multiplications, see Section 3.4, especially formula (3.20). We compute the attention score by

$$180 \quad (2.4) \quad \gamma(\mathbf{x}, \tilde{\mathbf{x}}, t; u) = \text{Softmax}_2 \left(\frac{1}{\sqrt{|\Omega_y|}} \left\langle Q(\mathbf{x}, \cdot, t; u), K(\tilde{\mathbf{x}}, \cdot, t; u) \right\rangle_{\Omega_y} \right),$$

181 where $\gamma(\mathbf{x}, \tilde{\mathbf{x}}, t)$ is defined on $\Omega_x \times \Omega_x \times [0, T]$, Softmax_2 represents the softmax function along
 182 the second dimension:

$$183 \quad \text{Softmax}_2(a(\mathbf{x}, \tilde{\mathbf{x}}, t)) = \frac{\exp(a(\mathbf{x}, \tilde{\mathbf{x}}, t))}{\int_{\Omega_x} \exp(a(\mathbf{x}, \eta, t)) d\eta}.$$

184 The output of the attention layer is then obtained by taking the inner product of the atten-
 185 tion score γ and the feature V : $\left\langle \gamma(\mathbf{x}, \cdot, t; u), V(\cdot, \mathbf{y}, t; u) \right\rangle_{\Omega_x}$. We emphasize that the first two
 186 integral transformations in (2.1)–(2.2) are used to generate the attention score $\gamma(\mathbf{x}, \tilde{\mathbf{x}}, t; u)$ and
 187 the last integral transformation (2.3) will be used for feature extractions. All three integral
 188 kernels will be learned during training, see explanations later in this section.

189 **2.3. Mathematical formulations for the Layer Normalization.** The layer normalization
 190 [1] in the Transformer [46] and also many other neural networks [60, 57, 51, 7] can be mathe-
 191 matically described as a projection of a function to a set with given mean value σ_1 and variance
 192 σ_2^2 . Define sets

$$193 \quad S_1(\sigma_1, \sigma_2) = \left\{ u : \frac{1}{|\Omega_y|} \int_{\Omega_y} u(\mathbf{x}, \xi, t) d\xi = \sigma_1, \right. \\
 194 \quad (2.5) \quad \left. \frac{1}{|\Omega_y|} \int_{\Omega_y} (u(\mathbf{x}, \xi, t) - \sigma_1)^2 d\xi = \sigma_2^2 \right\},$$

$$195 \quad (2.6) \quad S_2 = \{u : u \geq 0\},$$

196 and their corresponding indicator functions

$$197 \quad (2.7) \quad I_{S_1(\sigma_1, \sigma_2)}(u) = \begin{cases} 0 & \text{if } u \in S_1, \\ +\infty & \text{otherwise,} \end{cases} \quad I_{S_2}(u) = \begin{cases} 0 & \text{if } u \in S_2, \\ +\infty & \text{otherwise.} \end{cases}$$

198 Later, we shall show that layer normalization is just a projection of a function into the set
 199 $S_1(\sigma_1, \sigma_2)$, see Section 3.3.2 and the ReLU activation function is just a projection of a function
 200 to the set S_2 , see Section 3.3.3.

201 **2.4. The continuous Transformer.** In addition to the control variables $W^Q(\mathbf{y}, \tilde{\mathbf{y}}, t)$, $W^K(\mathbf{y}, \tilde{\mathbf{y}}, t)$,
 202 $W^V(\mathbf{y}, \tilde{\mathbf{y}}, t)$, let us also introduce control variables $\{W_j(\tilde{\mathbf{y}}, \mathbf{y}, t)\}_{j=1}^J$ on $\Omega_y \times \Omega_y \times [0, T]$, and
 203 $b(\mathbf{x}, t)$, $b_j(\mathbf{x}, t)$ on $\Omega_x \times [0, T]$. We call the following continuous integral equation a **Continuous**
 204 **Transformer**:

$$205 \quad (2.8) \quad \left\{ \begin{array}{l} u_t = \underbrace{\left\langle \gamma(\mathbf{x}, \cdot, t; u), V(\cdot, \mathbf{y}, t; u) \right\rangle_{\Omega_x}}_{\text{I: attention}} + \underbrace{\partial I_{S_1(\sigma_1(t), \sigma_2(t))}(u)}_{\text{II: layer normalization}} \\ \quad + \underbrace{\sum_{j=1}^J \left(\left\langle W_j(\cdot, \mathbf{y}, t), u(\mathbf{x}, \cdot, t) \right\rangle_{\Omega_y} + b_j(\mathbf{x}, t) \right)}_{\text{III: fully connected network}} + \partial I_{S_2}(u), \quad t \in (0, T], \\ u(\mathbf{x}, \mathbf{y}, 0) = f(\mathbf{x}, \mathbf{y}). \end{array} \right.$$

206 We shall show in the rest of this paper that the Transformer proposed in [46] is just a dis-
 207 cretization of this continuous Transformer. In (2.8), f is some initial state, T is a fixed time
 208 chosen by the user. Following conventions in the literature, we denote all the control variables
 209 as

$$210 \quad (2.9) \quad \theta = \{W^Q(\mathbf{x}, \mathbf{y}, t), W^K(\mathbf{x}, \mathbf{y}, t), W^V(\mathbf{x}, \mathbf{y}, t), \{W_j(\mathbf{x}, \mathbf{y}, t), b_j(\mathbf{x}, t)\}_{j=1}^J, \sigma_1(t), \sigma_2(t)\}.$$

211 The Continuous Transformer is the mapping:

$$212 \quad (2.10) \quad \mathcal{N}_\theta : f \mapsto u(\mathbf{x}, \mathbf{y}, T).$$

213 In Section 3, we will discretize (2.8) by the operator-splitting method for the time variable
 214 and use a uniform grid for the spatial variables \mathbf{x}, \mathbf{y} . After discretization, the continuous
 215 Transformer becomes the Transformer proposed in [46]. In relation to this equivalence, the
 216 operations in the right-hand side of (2.8) can be classified into three sets:

- 217 • Operations in Set I correspond to attentions in Transformers. It computes the attention
218 score $\gamma(\mathbf{x}, \tilde{\mathbf{x}}, t; u)$ and applies it to the extracted features in $V(\mathbf{x}, \mathbf{y}, t)$.
- 219 • Operations in Set II correspond to layer normalization. For each token location \mathbf{x} ,
220 it requires function u to have mean σ_1 and variance σ_2^2 along the \mathbf{y} direction. In
221 other words, for each token, the corresponding series is expected to have mean σ_1 and
222 variance σ_2^2 .
- 223 • Operations in Set III correspond to feedforward networks with J layers. The first com-
224 ponent corresponds to linear layers, each of which contains one linear transformation
225 W_j and a bias b_j . The second component corresponds to activation functions. This
226 term projects u to S_1 , which can be realized by ReLU functions.

227 **2.5. The continuous control problem for learning.** Our objective is to optimize θ so
228 that given an input $f(\mathbf{x}, \mathbf{y})$, the control problem will drive $u(\mathbf{x}, \mathbf{y}, T)$ to a desired state v .
229 Specifically, given a dataset $\{(u_i, v_i)\}_{i=1}^B$, where u_i is the input and v_i is the target state, let
230 $\ell(\cdot, \cdot)$ be a loss function measuring the discrepancy between its arguments. Let $\theta, \mathcal{N}_\theta$ be as
231 defined in (2.9)–(2.10), the training process with the loss function ℓ is in fact a process to
232 optimize θ for the following minimization problem:

$$233 \quad (2.11) \quad \min_{\theta} \frac{1}{B} \sum_{i=1}^B \ell(\mathcal{N}_\theta(u_i), v_i).$$

234 This is essentially a PDE-constrained optimization problem of the following form:

$$235 \quad (2.12) \quad \min_{\theta} \frac{1}{B} \sum_{i=1}^B \ell(w_i, v_i), \text{ under constraint } w_i = \mathcal{N}_\theta(u_i) \text{ solves the control problem (2.8).}$$

236 **3. Discretizations of the Continuous Problem.** Following the idea of PottsMGNet [42],
237 in this section, we use some proper operator splitting methods to discretize the time variable for
238 (2.8). For a review about operator splitting methods, please refer to [12, 13, 14]. Specifically,
239 we discretize the temporal domain of (2.8) by sequential splitting, c.f. [14, Sec. 2.2]. The
240 number of sequential splittings corresponds to the number of layers in Transformers. Spatial
241 discretization for \mathbf{x} and \mathbf{y} will be discussed in Section 3.4. The discretization of \mathbf{x} determines
242 the number of tokens in the input, and the discretization of \mathbf{y} determines the size of the
243 embedding vector for each token.

3.1. The overall time discretization. Let $\{t^n\}_{n=0}^{N_t}$ be the set of time grids with time step
 $\Delta t = T/N_t$. We denote the numerical solution of (2.11) at time t^n by u^n . We will design a
numerical scheme to approximate \mathcal{N}_θ . Specifically, for each t^n , we design a discrete propagator
 $\bar{\mathcal{N}}^n$ propagating u^{n-1} to u^n . Then, problem (2.8) is numerically approximated by

$$u(t^{N_t}) \approx u^{N_t} = \bar{\mathcal{N}}^{N_t} \circ \bar{\mathcal{N}}^{N_t-1} \circ \dots \circ \bar{\mathcal{N}}^1(f),$$

244 see Figure 1(b) for an illustration.

In each $\bar{\mathcal{N}}^n$, time-discretized control variables θ will be used. We denote the control
variables used in $\bar{\mathcal{N}}^n$ by θ^n . To show the dependence of $\bar{\mathcal{N}}^n$ on the control variable θ^n , we will

also use $\bar{\mathcal{N}}_{\theta^n}^n$ to denote the propagator $\bar{\mathcal{N}}^n$. From the following subsections, c.f. §3.4, we will see that

$$\theta^n = \{W^{Q,n-1}, W^{K,n-1}, W^{V,n-1}, \{W_j^{n-1}, b_j^{n-1}\}_{j=1}^J, \sigma_1^{n-1}, \sigma_2^{n-1}\}.$$

245 Denote all of the discretized variables by $\bar{\theta}$, i.e. $\bar{\theta} = \{\theta^n\}_{n=1}^{N_t}$. The whole evolution procedure
 246 from f to u^{N_t} by $\bar{\mathcal{N}}_{\bar{\theta}}$, i.e, $\bar{\mathcal{N}}_{\bar{\theta}} = \bar{\mathcal{N}}_{\theta^{N_t}}^{N_t} \circ \bar{\mathcal{N}}_{\theta^{N_t-1}}^{N_t-1} \circ \dots \circ \bar{\mathcal{N}}_{\theta^1}^1$. Given a dataset $\{(u_i, v_i)\}_{i=1}^B$, then
 247 the discrete training problem is to solve the following constrained optimization problem:

$$248 \quad (3.1) \quad \min_{\bar{\theta}} \frac{1}{B} \sum_{i=1}^B \ell(\bar{\mathcal{N}}_{\bar{\theta}}(u_i), v_i).$$

249 **3.2. Time Discretizations by Operator Splitting Method.** In the following sections, we
 250 illustrate our numerical scheme for propagator $\bar{\mathcal{N}}_{\theta^1}^1$. Propagator $\bar{\mathcal{N}}_{\theta^n}^n$ for $n = 2, \dots, N_t$ is exactly
 251 the same by replacing θ^1 by θ^n . Our numerical scheme is designed based on operator splitting
 252 methods. For simplicity, we take $\Delta t = 1$. It is often used in data analysis and image processing
 253 literature and is equivalent to changing the network parameters properly. In the following,
 254 with a slight abuse of notation, we denote the time-discretized version of any time-dependent
 255 function $a(\cdot, t)$ by $a^n(\cdot) := a(\cdot, t^n)$, where t^n is the n -th time step. Using the Lie scheme [14,
 256 Sec. 2.2] as explained in Appendix C, given $u^0 = f$, we compute u^1 by $M = 4 + J$ substeps:
 257 **Substep 1:** Solve $u^{1/M}$ from

$$258 \quad (3.2) \quad u^{1/M} - u^0 = \left\langle \gamma^0(\mathbf{x}, \cdot; u^0), V^0(\cdot, \mathbf{y}; u^0) \right\rangle_{\Omega_x}.$$

259 **Substep 2:** Compute $u^{2/M}$ from

$$260 \quad (3.3) \quad u^{2/M} - u^{1/M} = \partial I_{S_1(\sigma_1^0, \sigma_2^0)}(u^{2/M}).$$

261 **Substep 2+j :** Compute $u^{(2+j)/M}$ sequentially for $j = 1, \dots, J$ from

$$262 \quad (3.4) \quad u^{(2+j)/M} - u^{(1+j)/M} = \left\langle W_j^0(\cdot, \mathbf{y}), u^{(1+j)/M}(\mathbf{x}, \cdot) \right\rangle_{\Omega_y} + b_j^0(\mathbf{x}) + \partial I_{S_2}(u^{(2+j)/M}).$$

263 **Substep 3+J:** Compute $u^{(3+J)/M}$ from

$$264 \quad (3.5) \quad u^{(3+J)/M} = \frac{1}{2}(u^{(2+J)/M} + u^{2/M})$$

265 **Substep 4+J:** Compute u^1 from

$$266 \quad (3.6) \quad u^1 - u^{(3+J)/M} = \partial I_{S_1(\sigma_1^0, \sigma_2^0)}(u^1)$$

267 Note that in (3.2)–(3.4) and (3.6), the left hand side $u^{(k+1)/M} - u^{k/M} = \frac{u^{(k+1)/M} - u^{k/M}}{\Delta t}$ with
 268 $\Delta t = 1$ is the finite difference approximation of u_t .

269 Substep 1 to Substep (4+J) evolve u^0 to u^1 by sequentially passing it through all operators
 270 in (2.8), corresponding to all components in Transformer from input to output. Substep 1
 271 corresponds to the attention layer and the skip connection. Substep 2 corresponds to layer

272 normalization. Substep 3 to Substep $(4 + J)$ correspond to the J layers in the fully connected
 273 network. Substep $3 + J$ corresponds to the skip connection after the fully connected network.
 274 The last substep corresponds to the final layer normalization. From here, we can already see
 275 that each substep in the Lie splitting scheme corresponds to one layer in the Transformer. In
 276 [46], it used $J = 2$ and thus $M = 6$. The different layers in the Transformer is just a sequential
 277 computation of the substep solutions:

$$\begin{aligned}
 278 \quad (3.7) \quad u^0 \quad (\text{input token}) &\mapsto u^{1/6}(\text{attention}) \mapsto u^{2/6}(\text{normalization layer}) \\
 279 &\mapsto u^{3/6} \text{ and } u^{4/6}(\text{fully connected layer with ReLU}) \\
 280 &\mapsto u^{5/6}(\text{skip connection layer}) \mapsto u^1(\text{normalization layer}).
 \end{aligned}$$

281 **3.3. Solutions to Each Subproblem.** In the following, we supply the details for the solu-
 282 tion to each subproblem in the splitting scheme (3.2)–(3.6).

283 **3.3.1. Solution to Subproblem (3.2).** We have

$$284 \quad (3.8) \quad u^{1/M} = u^0 + \int_{\Omega_x} \text{Softmax}_2 \left(\frac{1}{\sqrt{|\Omega_y|}} \left\langle Q^0(\mathbf{x}, \cdot; u^0), K^0(\eta, \cdot; u^0; \cdot) \right\rangle_{\Omega_y} \right) V^0(\eta, \mathbf{y}; u^0) d\eta.$$

285 This step is a continuous version of the attention layer in [46]. It looks very simple and uses
 286 two integral transformations to get the attention score through a softmax operator. This score
 287 function has values between $[0, 1]$ and it is multiplied with the $V^0(\eta, \mathbf{y}; u^0)$, which are the
 288 extracted features from the input data through another integral transformation.

289 **3.3.2. Solution to Subproblem (3.3) and (3.6).** For $u^{2/M}$ and u^1 in problem (3.3) and
 290 (3.6), they are in the form of

$$291 \quad (3.9) \quad u - v = \partial I_{S_1(\sigma_1, \sigma_2)}(u)$$

292 for the given function $v(\mathbf{x}, \mathbf{y})$ and constants σ_1, σ_2 . It is easy to derive that the solution u of
 293 (3.9) solves the following problem:

$$294 \quad (3.10) \quad u = \underset{\bar{u} \in S_1(\sigma_1, \sigma_2)}{\operatorname{argmin}} \frac{1}{2} \|\bar{u} - v\|_{\Omega_y}^2.$$

295 This shows that u is a projection of v to the set $S_1(\sigma_1, \sigma_2)$. The following theorem gives a
 296 closed-form solution formula for u :

297 **Theorem 3.1.** *The solution to problem (3.10) is given as:*

$$298 \quad (3.11) \quad u(\mathbf{x}, \mathbf{y}) = \frac{v(\mathbf{x}, \mathbf{y}) - \alpha(\mathbf{x})}{\sqrt{\beta(\mathbf{x})}} \sigma_2 + \sigma_1,$$

$$299 \quad (3.12) \quad \text{with } \alpha(\mathbf{x}; v) = \frac{1}{|\Omega_y|} \int_{\Omega_y} v(\mathbf{x}, \xi) d\xi, \quad \beta(\mathbf{x}; v) = \frac{1}{|\Omega_y|} \int_{\Omega_y} (v(\mathbf{x}, \xi) - \alpha(\mathbf{x}))^2 d\xi.$$

300 Theorem 3.1 is proved in Appendix B. We remark that the expression (3.12) recovers the
 301 normalization layer in the continuous setting. According to Theorem 3.1, we see that the
 302 solutions $u^{2/M}$ and $u^{(4+J)/M}$ are explicitly given as:

$$303 \quad (3.13) \quad u^{2/M}(\mathbf{x}, \mathbf{y}) = \frac{u^{1/M}(\mathbf{x}, \mathbf{y}) - \alpha(\mathbf{x}; u^{1/M})}{\sqrt{\beta(\mathbf{x}; u^{1/M})}} \sigma_2^0 + \sigma_1^0,$$

$$304 \quad (3.14) \quad u^1(\mathbf{x}, \mathbf{y}) = \frac{u^{(3+J)/M}(\mathbf{x}, \mathbf{y}) - \alpha(\mathbf{x}; u^{(3+J)/M})}{\sqrt{\beta(\mathbf{x}; u^{(3+J)/M})}} \sigma_2^0 + \sigma_1^0.$$

305 **3.3.3. Solution to Subproblem (3.4).** Problem (3.4) is a semi-implicit equation for $u^{(2+j)/M}$.
 306 It can be solved exactly by the following sequential splitting

$$307 \quad (3.15) \quad \begin{cases} \bar{u}^{(2+j)/M} = u^{(1+j)/M} + \left\langle W_j^0(\cdot, \mathbf{y}), u^{(1+j)/M}(\mathbf{x}, \cdot) \right\rangle_{\Omega_y} + b_j^0(\mathbf{x}), \\ u^{(2+j)/M} - \bar{u}^{(2+j)/M} = \partial I_{S_2}(u^{(2+j)/M}). \end{cases}$$

308 The first equation in (3.15) is linear in $\bar{u}^{(2+j)/M}$, corresponding to a linear layer in a feed-
 309 forward network in continuous settings. The second equation solves

$$310 \quad (3.16) \quad u^{(2+j)/M} = \operatorname{argmin}_{v \in S_1} \int_{\Omega_x} \int_{\Omega_y} |v - \bar{u}^{(2+j)/M}|^2 d\mathbf{y} d\mathbf{x}.$$

311 It is easy to see that its solution can be computed point-wisely:

$$312 \quad (3.17) \quad u^{(2+j)/M} = \max\{\bar{u}^{(2+j)/M}, 0\} = \operatorname{ReLU}(\bar{u}^{(2+j)/M}),$$

313 which corresponds ReLU activation. As a result, Substep $2+j$ for $j = 1, \dots, J$ is a feedforward
 314 network layer activated by ReLU.

315 **3.4. Spatial Discretizations.** In this subsection, we discretize the solution discussed in
 316 Section 3.3 spatially. We will show that after spatial discretizations, the splitting scheme
 317 (3.2)–(3.6) exactly recovers the Transformer architecture in [46].

Suppose $\Omega_x = [0, L_x]$, $\Omega_y = [0, L_y]$. We uniformly discretize Ω_x, Ω_y by N_x and N_y grids in the \mathbf{x} and \mathbf{y} coordinates, respectively. For simplicity, we suppose $L_x = N_x, L_y = N_y$, leading to discretization steps $\Delta x = L_x/N_x = 1$, $\Delta y = L_y/N_y = 1$, and the grids by $\mathbf{x}_k = k, \mathbf{y}_k = k$. We use $\mathbf{1}_m$ to denote an m -dimensional vector with elements 1. For a matrix $\mathbf{W} \in \mathbb{R}^{N_x \times N_y}$ and a vector $\mathbf{b} \in \mathbb{R}^{N_x}$, the matrix $\mathbf{z} = (\mathbf{W} + \mathbf{b}) \in \mathbb{R}^{N_x \times N_y}$ is computed so that

$$\mathbf{z}_{k,l} = W_{k,l} + b_k.$$

318 This is an easy notation to add a vector to a matrix in the discrete setting. With the spatial
 319 discretizations above, we define the discrete integrals (sum) as

$$320 \quad (3.18) \quad \int_{\bar{\Omega}_x} f d\mathbf{x} = \sum_{k=1}^{N_x} f_k, \quad \int_{\bar{\Omega}_x} \int_{\bar{\Omega}_y} f(\mathbf{x}, \mathbf{y}) d\mathbf{x} d\mathbf{y} = \sum_{k=1}^{N_x} \sum_{l=1}^{N_y} f_{k,l}.$$

321 For any functions $v(\mathbf{x}), w(\mathbf{x})$ defined on Ω_x , denote their discrete counterpart by \mathbf{v}, \mathbf{w} . We
 322 denote the discretized inner product between \mathbf{v} and \mathbf{w} as

$$323 \quad (3.19) \quad \langle \mathbf{v}, \mathbf{w} \rangle_{\bar{\Omega}_x} = \mathbf{v}^\top \mathbf{w} = \sum_{k=1}^{N_x} v_k w_k.$$

324 The discretized integrals and inner products $\langle \cdot, \cdot \rangle_{\bar{\Omega}_y}$ and $\langle \cdot, \cdot \rangle_{\bar{\Omega}_x \times \bar{\Omega}_y}$ are defined similarly.

325 **3.4.1. Computing the Discrete Solutions $u^{1/M}$ in (3.8).** Before we discretize the formula
 326 (3.8), we first discretize some related quantities and operations. To compute the attention
 327 scores, we discretize the integral transformations as follows:

$$328 \quad (3.20) \quad \mathbf{Q}^0(\mathbf{u}^0) = \mathbf{u}^0 \mathbf{W}^{Q,0}, \quad \mathbf{K}^0(\mathbf{u}^0) = \mathbf{u}^0 \mathbf{W}^{K,0}, \quad \mathbf{V}^0(\mathbf{u}^0) = \mathbf{u}^0 \mathbf{W}^{V,0}.$$

329 Here, \mathbf{u}^0 is in $\mathbb{R}^{N_x \times N_y}$ and $\mathbf{W}^{Q,0}, \mathbf{W}^{K,0}, \mathbf{W}^{V,0}$ are in $\mathbb{R}^{N_y \times N_y}$, and the expressions such as
 330 $\mathbf{u}^0 \mathbf{W}^{Q,0}$ represent standard matrix multiplications that yield new matrices in $\mathbb{R}^{N_x \times N_y}$. We
 331 emphasize that the multiplication used here is the conventional matrix product. For $\mathbf{A} =$
 332 $\{a_{k,l}\}_{k,l} \in \mathbb{R}^{N_x \times N_x}$, the discretized version of $\text{Softmax}_2(\mathbf{A})$ is given as

$$333 \quad (3.21) \quad (\text{Softmax}_{2,\text{dis}}(\mathbf{A}))_{k,l} = \frac{\exp(a_{k,l})}{\sum_{l=1}^{N_x} \exp(a_{k,l})}.$$

334 Then the updating formula (3.8) for $u^{1/M}$ in the discrete setting is

$$335 \quad (3.22) \quad \mathbf{u}^{1/M} = \mathbf{u}^0 + \text{Softmax}_{2,\text{dis}} \left(\mathbf{Q}^0(\mathbf{u}^0) (\mathbf{K}^0(\mathbf{u}^0))^\top \right) \mathbf{V}^0(\mathbf{u}^0).$$

336 We emphasize again that this is exactly the attention layer in [46]. This is just a compact
 337 mathematical expression for it. All the matrices $\mathbf{W}^{Q,0}, \mathbf{W}^{K,0}, \mathbf{W}^{V,0}$ will be learned in the
 338 training process.

339 **3.4.2. Computing the Discrete Solutions $u^{2/M}$ and $u^{(4+J)/M}$ in (3.9).** In Theorem 3.1,
 340 for any $\mathbf{v} = \{v_{k,l}\}_{k,l} \in \mathbb{R}^{N_x \times N_y}$, we discretize $\alpha(\mathbf{x})$ and $\beta(\mathbf{x})$ as

$$341 \quad (3.23) \quad \alpha_k(\mathbf{v}) = \frac{1}{L_y} \sum_{l=1}^{N_y} v_{k,l}, \quad \beta_k(\mathbf{v}) = \frac{1}{L_y} \sum_l (v_{k,l} - \alpha_k(\mathbf{v}))^2$$

342 for $k = 1, \dots, N_x$. Then we compute the discrete solutions $u^{2/M}$ and $u^{(4+J)/M}$ as

$$343 \quad (3.24) \quad u_{k,l}^{2/M} = \sigma_2^0 \frac{u_{k,l}^{1/M} - \alpha_k(\mathbf{u}^{1/M})}{\sqrt{\beta_k(\mathbf{u}^{1/M})}} + \sigma_1^0,$$

$$344 \quad (3.25) \quad u_{k,l}^{n+(4+J)/M} = \sigma_2^0 \frac{u_{k,l}^{(3+J)/M} - \alpha_k(\mathbf{u}^{(3+J)/M})}{\sqrt{\beta_k(\mathbf{u}^{(3+J)/M})}} + \sigma_1^0,$$

345 for $k = 1, \dots, N_x$ and $l = 1, \dots, N_y$.

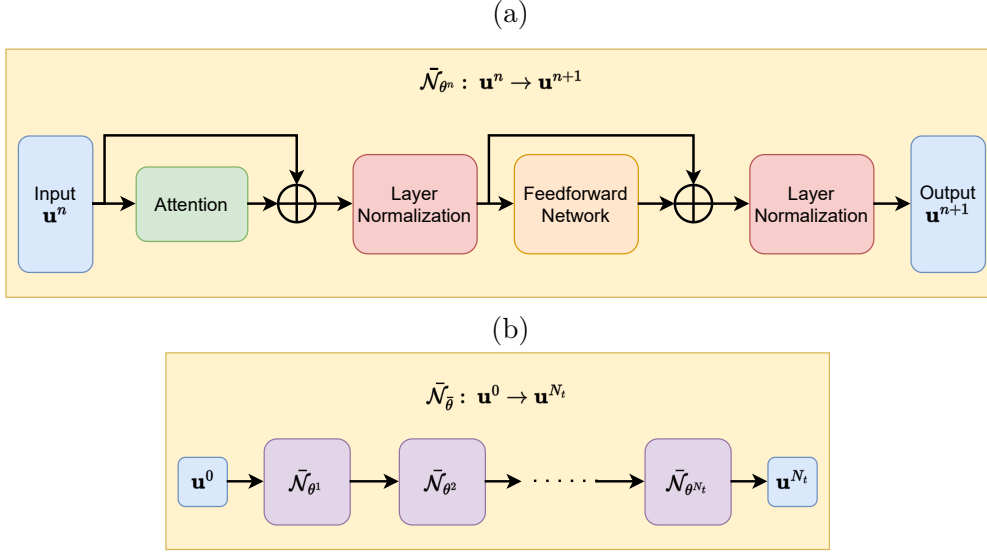


Figure 1. (a) Illustration of the discretized operator-splitting scheme (3.2)–(3.6). (b) Illustration of the whole evolution.

346 **3.4.3. Computing the Discrete Solutions $u^{(3+J)/M}$ in (3.5).** For $u^{(3+J)/M}$, its discrete
347 analogue is computed as

$$348 \quad (3.26) \quad \mathbf{u}^{(3+J)/M} = \frac{1}{2}(\mathbf{u}^{(2+J)/M} + \mathbf{u}^{2/M}).$$

349 **3.4.4. Computing the Discrete Solutions $u^{(2+j)/M}$ in (3.15).** In (3.15), we first compute

$$350 \quad (3.27) \quad \bar{\mathbf{u}}^{(2+j)/M} = \mathbf{u}^{(1+j)/M} + \mathbf{u}^{(1+j)/M} \mathbf{W}_j^0 + \mathbf{b}_j^0,$$

351 and then update

$$352 \quad (3.28) \quad \mathbf{u}^{(2+j)/M} = \text{ReLU}(\bar{\mathbf{u}}^{(2+j)/M}).$$

353 The structure of the discretized scheme (3.2)–(3.6) is illustrated in Figure 1(a), and it exactly
354 recovers the well-known Transformer architecture introduced in [46], see the next section for
355 detailed explanations.

356 4. Mathematical Explanations of Transformers.

357 **4.1. Connections to Transformer Encoder in [46].** Transformer architecture was first
358 proposed in [46] for natural language processing. After time and spatial discretization, our
359 splitting scheme recovers the well-known Transformer with single-head attention in [46]. We
360 supply the details of the explanations in the following. Let us first show that one iteration
361 $\mathbf{u}^0 \rightarrow \mathbf{u}^1$ realizes one Transformer block:

- 362 • Input \mathbf{u}^0 : With our discretization, $\mathbf{u}^0 \in \mathbb{R}^{N_x \times N_y}$. It is an input with N_x tokens. Each
363 token is represented as an integer-vector with N_y numbers.

- 364 • Scaled dot-product attention: Equation (3.22) is a single head attention with a skip
365 connection. The second term in the right-hand side of (3.22) is a single-head attention
366 with query matrix \mathbf{Q}^0 , key matrix \mathbf{K}^0 and value matrix \mathbf{V}^0 , whose weight matrices are
367 $\mathbf{W}^{Q,0}$, $\mathbf{W}^{K,0}$ and $\mathbf{W}^{V,0}$, respectively. It recovers the single-head attention in [46]. The
368 first term in the right-hand side is a skip connection, adding together the input and
369 output of attention.
- 370 • Layer normalization: Equation (3.24)–(3.25) perform layer normalization. For each
371 token, (3.24)–(3.25) normalize the input along the embedding direction to have mean
372 σ_1 and variance σ_2^2 .
- 373 • Position-wise feedforward network: Equation (3.27)–(3.28) realizes a feedforward net-
374 work with J layers activated by ReLU. In particular, (3.27) can be rewritten as

$$375 \quad \bar{\mathbf{u}}^{(2+j)/M} = \mathbf{u}^{(1+j)/M}(\mathbf{I} + \mathbf{W}_j^0) + \mathbf{b}_j^0,$$

376 where \mathbf{I} denotes the identity matrix. It is a linear layer with weight matrix $(\mathbf{I} + \mathbf{W}_j^0)$
377 and bias \mathbf{b}_j^0 . Set $J = 2$. Equation (3.27)–(3.28) recovers the feedforward network in
378 [46].

- 379 • Skip connection: Equation (3.26) is a relaxation step in the splitting scheme. It realizes
380 the skip connection by averaging the input and output of the feedforward network.

381 The discussion above shows that one time step of the splitting scheme supplied in Section
382 3.2 realizes the Transformer of [46]. In Section 3.2, the time domain is discretized into N_t
383 time steps. Thus, the whole operator-splitting scheme is equivalent to compositions of N_t
384 Transformer blocks. Since the learnable variables $\theta^n = \{\mathbf{W}^{Q,n}, \mathbf{W}^{K,n}, \mathbf{W}^{V,n}, \{\mathbf{W}_j^n, \mathbf{b}_j^n\}_{j=1}^J\}$
385 depends on time, they are different for different time steps. Setting $N_t = 6$ recovers the
386 architecture in [46]. Furthermore, solving (3.1) for the control variables is equivalent to training
387 the network.

388 **4.2. Connections to ViT in [8].** Based on the Transformer encoder in [46], vision Trans-
389 former (ViT) was proposed in [8] for image classification. In ViT, each image is cropped
390 into patches, each of which is taken as a token. ViT has a similar architecture as the Trans-
391 former encoder in [46], except that it has an additional input embedding layer with learnable
392 parameters, and a linear layer at the end for output.

393 In order to use (3.2)–(3.6) to provide a mathematical explanation of ViT, we need to
394 incorporate the additional embedding layer and last linear layer with our algorithm. It can be
395 easily done by introducing data-driven data pre-processing and post-processing steps. For the
396 additional input embedding, we design a pre-processing step to generate the initial condition
397 f . Consider the discrete setting. Denote $\mathbf{R} \in \mathbb{R}^{(N_x-1) \times D}$ as the collection of $(N_x - 1)$ patches
398 of an image, where each row of \mathbf{R} corresponds to one flattened patch of dimension D . Let
399 $\mathbf{E} \in \mathbb{R}^{D \times N_y}$ be an embedding matrix and $\mathbf{v}_0 \in \mathbb{R}^{1 \times N_y}$ be a vector where both are learnable.
400 The vector \mathbf{v}_0 represents the extra learnable embedding. We generate \mathbf{u}_0 as

$$401 \quad (4.1) \quad \mathbf{u}_0 = F(\mathbf{R}) = \begin{bmatrix} \mathbf{v}_0 \\ \mathbf{R}\mathbf{E} \end{bmatrix}.$$

402 The last linear layer in ViT can be taken as a data post-processing. For any $\mathbf{u} \in \mathbb{R}^{N_x \times N_y}$ and

403 weight matrix $\mathbf{B} \in \mathbb{R}^{N_y \times d}$. Define

$$404 \quad (4.2) \quad G(\mathbf{B}; \mathbf{u}) = [1 \quad \mathbf{0}_{1 \times (N_x - 1)}] \mathbf{u} \mathbf{B},$$

where $\mathbf{0}_{1 \times (N_x - 1)}$ is a zero vector of size $(N_x - 1)$. Let $\bar{\theta}, \bar{\mathcal{N}}$ be defined as in Section 3.1. We denote

$$\tilde{\theta} = \{\bar{\theta}, \mathbf{v}_0, \mathbf{E}, \mathbf{B}\}$$

405 as the collection of all learnable parameters. With the data pre- and post-processing, we define
406 the whole operation

$$407 \quad (4.3) \quad \bar{\mathcal{N}}_{\tilde{\theta}}^{\text{ViT}} = G \circ \bar{\mathcal{N}} \circ F.$$

408 One can readily check that $\bar{\mathcal{N}}_{\tilde{\theta}}^{\text{ViT}}$ recovers the ViT in [8]. Given a dataset $\{(\mathbf{u}_i, \mathbf{v}_i)\}_{i=1}^B$, we
409 learn all parameters by solving

$$410 \quad (4.4) \quad \min_{\tilde{\theta}} \frac{1}{B} \sum_{i=1}^B \ell(\bar{\mathcal{N}}_{\tilde{\theta}}^{\text{ViT}}(\mathbf{u}_i), \mathbf{v}_i),$$

411 which is equivalent to training a ViT.

412 **5. Multi-Head Attention.** As discussed in Section 4, the proposed splitting scheme for
413 problem (2.8) is the Transformer architecture with one-head attention. In this section, we
414 extend (2.8) and splitting scheme (3.2)–(3.6) so that it realizes multi-head attention.

415 **5.1. Control Problem with Continuous Head Dimension.** We first consider a continuous-
416 head setting. After discretizations, we will show that it recovers multi-head attention. In the
417 continuous-head setting, we take head as an additional dimension. To accommodate the new di-
418 mension, we update the definition of W^Q, W^K, W^V to $W^Q(\mathbf{y}, \tilde{\mathbf{y}}, h, t), W^K(\mathbf{y}, \tilde{\mathbf{y}}, h, t), W^V(\mathbf{y}, \tilde{\mathbf{y}}, h, t)$ ■
419 on $\Omega_y \times \Omega_y \times \Omega_h \times [0, +\infty)$, where h is the variable for the head dimension with domain Ω_h .
420 Similar to (2.1)–(2.3), we define

$$\begin{aligned} 421 \quad Q(\mathbf{x}, \mathbf{y}, h, t; u) &= \left\langle W^Q(\cdot, \mathbf{y}, h, t), u(\mathbf{x}, \cdot, t) \right\rangle_{\Omega_y}, \\ 422 \quad K(\mathbf{x}, \mathbf{y}, h, t; u) &= \left\langle W^K(\cdot, \mathbf{y}, h, t), u(\mathbf{x}, \cdot, t) \right\rangle_{\Omega_y}, \\ 423 \quad (5.1) \quad V(\mathbf{x}, \mathbf{y}, h, t; u) &= \left\langle W^V(\cdot, \mathbf{y}, h, t), u(\mathbf{x}, \cdot, t) \right\rangle_{\Omega_y}. \end{aligned}$$

424 We update $\gamma(\mathbf{x}, \tilde{\mathbf{x}}, t; u)$ in (2.4) to

$$425 \quad (5.2) \quad \gamma(\mathbf{x}, \tilde{\mathbf{x}}, h, t; u) = \text{Softmax}_2 \left(\frac{1}{\sqrt{|\Omega_y|}} \left\langle Q(\mathbf{x}, \cdot, h, t; u), K(\tilde{\mathbf{x}}, \cdot, h, t; u) \right\rangle_{\Omega_y} \right).$$

426 The control problem (2.8) is updated to

$$\begin{cases}
 427 \quad (5.3) & \left\{ \begin{aligned}
 & u_t = \underbrace{\int_{\Omega_h} \langle \gamma(\mathbf{x}, \cdot, h, t; u), V(\cdot, \mathbf{y}, h, t; u) \rangle_{\Omega_x} dh}_{\text{I-a: attention}} + \underbrace{\partial I_{S_1(\sigma_1, \sigma_2)}(u)}_{\text{II: layer normalization}} \\
 & + \underbrace{\sum_{j=1}^J \left(\langle W_j(\cdot, \mathbf{y}, t), u(\mathbf{x}, \cdot, t) \rangle_{\Omega_y} + b_j(\mathbf{x}, t) \right)}_{\text{III: fully connected network}} + \partial I_{S_2}(u) \text{ for } t \in (0, T], \\
 & u(\mathbf{x}, \mathbf{y}, 0) = f(\mathbf{x}, \mathbf{y}).
 \end{aligned} \right.
 \end{cases}$$

428 Compared to (2.8), problem (5.3) replace I by I-a, which introduces an additional integral with
 429 respect to h .

430 **5.2. Time Discretizations Through Operator-Splitting Method.** The time discretiza-
 431 tions and operator-splitting scheme for problem (5.3) are the same as those for problem (2.8)
 432 except that we update Substep 1 to

$$433 \quad (5.4) \quad u^{1/M} - u^0 = \int_{\Omega_h} \langle \gamma^0(\mathbf{x}, \cdot, h; u^0), V^0(\cdot, \mathbf{y}, h; u^0) \rangle_{\Omega_x} dh,$$

434 whose solution is given as

$$435 \quad u^{1/M} =$$

$$436 \quad (5.5) \quad u^0 + \int_{\Omega_h} \int_{\Omega_x} \text{Softmax}_2 \left(\frac{1}{\sqrt{|\Omega_y|}} \langle Q^0(\mathbf{x}, \cdot, h; u^0), K^0(\eta, \cdot, h; u^0) \rangle_{\Omega_y} \right) V^0(\eta, \mathbf{y}, h; u^0) d\eta dh.$$

437 **5.3. Discrete Solutions of $u^{1/M}$ in (5.5).** Suppose $\Omega_h = [0, L_h]$. In addition to the
 438 spatial discretization discussed in Section 3.4, we discretize the head domain into N_h grids.
 439 For simplicity, we suppose $L_h = N_h$.

440 For function $W^{K,0}(\mathbf{y}, \tilde{\mathbf{y}}, h)$, we denote $W_{k,l,m}^{Q,0} = W^{Q,0}(k, l, m)$, $\mathbf{W}_m^{Q,0} = \{W_{k,l,m}^{Q,0}\}_{k,l} \in$
 441 $\mathbb{R}^{N_y \times N_y}$ and $\mathbf{W}^{Q,0} = \{\mathbf{W}_m^{Q,0}\}_m \in \mathbb{R}^{N_y \times N_y \times N_h}$. We define $\mathbf{W}^{K,0}, \mathbf{W}^{V,0}$ similarly. We de-
 442 note the discretized integral transformation in the updated attention score as

$$443 \quad \mathbf{Q}_m^0(\mathbf{u}^0) = \mathbf{u}^0 \mathbf{W}_m^{Q,0}, \quad \mathbf{K}_m^0(\mathbf{u}^0) = \mathbf{u}^0 \mathbf{W}_m^{K,0}, \quad \mathbf{V}_m^0(\mathbf{u}^0) = \mathbf{u}^0 \mathbf{W}_m^{V,0},$$

444 and denote $\mathbf{Q}^0(\mathbf{u}^0) = \{\mathbf{Q}_m^0(\mathbf{u}^0)\}_m$, $\mathbf{K}^0(\mathbf{u}^0) = \{\mathbf{K}_m^0(\mathbf{u}^0)\}_m$, $\mathbf{V}^0(\mathbf{u}^0) = \{\mathbf{V}_m^0(\mathbf{u}^0)\}_m \in \mathbb{R}^{N_x \times N_y \times N_h}$. ■

445 The inner product $\langle Q^0(\mathbf{x}, \cdot, h; u^0), K^0(\tilde{\mathbf{x}}, \cdot, h; u^0) \rangle_{\Omega_y}$ is discretized as

$$446 \quad (5.6) \quad \langle \mathbf{Q}_m^0(\mathbf{u}^0), \mathbf{K}_m^0(\mathbf{u}^0) \rangle_{\tilde{\Omega}_y} = \mathbf{Q}_m^0(\mathbf{u}^0) (\mathbf{K}_m^0(\mathbf{u}^0))^\top.$$

447 Then the updating formula (5.5) for $u^{1/M}$ is discretized as

$$448 \quad (5.7) \quad \mathbf{u}^{1/M} = \mathbf{u}^0 + \sum_{m=1}^{N_h} \text{Softmax}_{2,\text{dis}} \left(\frac{1}{\sqrt{L_x}} \mathbf{Q}_m^0(\mathbf{u}^0) (\mathbf{K}_m^0(\mathbf{u}^0))^\top \right) \mathbf{V}_m^0(\mathbf{u}^0).$$

449 **5.4. Connections to Transformers.** Similar to the discussion in Section 4, with (5.7), the
 450 operator-splitting scheme for the control problem (5.3) realizes N_t encoder time steps in [46]
 451 with multi-head attention. To show this, we only need to demonstrate the equivalence between
 452 (5.7) and the multi-head attention.

453 The second term in the right-hand side of (5.7) sums over the embedding dimension and
 454 head dimension. It is a multi-head attention with query matrix \mathbf{Q}^0 , key matrix \mathbf{K}^0 and value
 455 matrix \mathbf{V}^0 , whose weight matrices are $\mathbf{W}^{Q,0}$, $\mathbf{W}^{K,0}$ and $\mathbf{W}^{V,0}$, respectively.

456 **6. Convolution Transformers for Video and Image Problems.** Image and video data
 457 inherently possess rich local spatial structures, which have been effectively exploited by con-
 458 volutional neural networks (CNNs) [37]. Convolutions serve as a powerful mechanism for
 459 capturing spatially localized patterns, making them especially well-suited for data defined on
 460 Cartesian grids.

461 In contrast, Transformer architectures were originally developed for language processing,
 462 where the data—represented by functions such as $u(\mathbf{x}, \mathbf{y}, t)$ —often lack explicit spatial struc-
 463 ture, as in the case of sentences or token sequences. For such unstructured data, attention-
 464 based mechanisms leveraging general integral transforms offer a flexible and effective modeling
 465 tool.

466 However, when applying Transformers to structured data such as images, videos, or other
 467 signals on regular grids, it has been shown that convolutional operations are significantly more
 468 efficient for feature extraction. This is due to their locality, weight sharing, and computational
 469 efficiency. Notably, convolutions can be interpreted as a special class of integral transformations
 470 with translation-invariant kernels and localized support.

471 Thanks to the integral formulation of the Q , K , and V operators in our framework (Sec-
 472 tions 2–3), it becomes straightforward to specialize these operators to convolutions. This al-
 473 lows our model to seamlessly incorporate convolutional structures, thereby leveraging domain-
 474 specific inductive biases and enhancing computational performance when applied to grid-
 475 structured data.

476 **6.1. Convolution Transformers.** We discuss the case for grayscale images. The framework
 477 can be extended to color images and videos by introducing additional variables to u .

478 We consider functions of interest $u(x, \mathbf{y}, t)$ where $x \in \Omega_x \subset \mathbb{R}$ is the index of tokens (image
 479 patches), $\mathbf{y} \in \Omega_y \subset \mathbb{R}^2$ are variables for spatial domain (indices of pixel location), and t is the
 480 time. For each token and any fixed time t , u is a two-dimensional image.

481 In this setting, we define $W^Q(\mathbf{y}, t), W^K(\mathbf{y}, t), W^V(\mathbf{y}, t) \in \Omega_y \times [0, +\infty)$ as convolution
 482 kernels. Operations Q, K, V are defined by convolution as

$$\begin{aligned}
 483 \quad Q(x, \mathbf{y}, t; u) &= W^Q(\cdot, t) * u(x, \cdot, t) = \int_{\Omega_y} W^Q(\boldsymbol{\xi}, t) u(x, \mathbf{y} - \boldsymbol{\xi}, t) d\boldsymbol{\xi}, \\
 484 \quad K(x, \mathbf{y}, t; u) &= W^K(\cdot, t) * u(x, \cdot, t) = \int_{\Omega_y} W^K(\boldsymbol{\xi}, t) u(x, \mathbf{y} - \boldsymbol{\xi}, t) d\boldsymbol{\xi}, \\
 485 \quad (6.1) \quad V(x, \mathbf{y}, t; u) &= W^V(\cdot, t) * u(x, \cdot, t) = \int_{\Omega_y} W^V(\boldsymbol{\xi}, t) u(x, \mathbf{y} - \boldsymbol{\xi}, t) d\boldsymbol{\xi}.
 \end{aligned}$$

486 The score function is computed as

$$487 \quad (6.2) \quad \gamma(x, \tilde{x}, t; u) = \text{Softmax}_2 \left(\frac{1}{\sqrt{|\Omega_y|}} \int_{\Omega_y} Q(x, \boldsymbol{\xi}, t; u) K(\tilde{x}, \boldsymbol{\xi}, t; u) d\boldsymbol{\xi} \right),$$

488 and the output of the attention layer is then given by: $\langle \gamma(x, \cdot, t; u), V(\cdot, \mathbf{y}, t; u) \rangle_{\Omega_x}$.

489 With the modification above, scheme (3.2)–(3.4) gives rise to convolutional Transformer.
490 Solvers discussed in Section 3.3 can still be used to solve subproblems in the new scheme.

491 **6.2. Connections to CvT in [50].** Convolutional vision Transformer (CvT) [50] incorpo-
492 rates Vit with convolutional neural networks (CNN), which utilizes both image local spatial
493 structures (by CNN) and global contexture (by Transformer). CvT consists of two compo-
494 nents, convolutional token embedding and convolutional Transformer block. The convolution
495 Transformer discussed in Section 6.1 realizes a one-head CvT with one stage and without
496 convolutional token embedding.

497 In fact, the convolutional token embedding is an additional convolution layer. Our operator-
498 splitting scheme can be easily revised to accomodate this layer. Let $W^C(\mathbf{x}, \mathbf{y}, t)$ be a convolu-
499 tion kernel and $b^C(\mathbf{x}, t)$ be a bias term. We consider the following control problem

$$500 \quad (6.3) \quad \left\{ \begin{array}{l} u_t = \underbrace{\langle \gamma(\mathbf{x}, \cdot, t; u), V(\cdot, \mathbf{y}, t; u) \rangle_{\Omega_x}}_{\text{I: attention}} + \underbrace{\partial I_{S_1(\sigma_1(t), \sigma_2(t))}(u)}_{\text{II: layer norml.}} + \underbrace{W^C(\cdot, t) * u(\mathbf{x}, \cdot, t) + b^C(\mathbf{x}, t)}_{\text{IV: conv. token embedding}} \\ \quad + \underbrace{\sum_{j=1}^J \left(\langle W_j(\cdot, \mathbf{y}, t), u(\mathbf{x}, \cdot, t) \rangle_{\Omega_y} + b_j(\mathbf{x}, t) \right)}_{\text{III: fully connected network}} + \partial I_{S_2}(u), \quad t \in (0, T], \\ u(\mathbf{x}, \mathbf{y}, 0) = f(\mathbf{x}, \mathbf{y}), \end{array} \right.$$

501 where the score function γ is computed as discussed in Section 6.1. For the operator-splitting
502 scheme (3.2)–(3.4), we add two additional substeps before (3.2). Given u^0 , we compute u^1 by
503 $M = 6 + J$ substeps:

504 **Substep 1:** Solve $\tilde{u}^{1/M}$ from

$$505 \quad (6.4) \quad u^{1/M} - u^0 = W^{C,0}(\cdot) * u^0(\mathbf{x}, \cdot) + b^{C,0}(\mathbf{x}).$$

506 **Substep 2:** Compute $u^{2/M}$ from

$$507 \quad (6.5) \quad u^{2/M} - u^{1/M} = \partial I_{S_1(\sigma_1^0, \sigma_2^0)}(u^{2/M}).$$

508 The remaining intermediate variables $u^{3/M}, \dots, u^1$ can be computed via (3.2)–(3.4), respec-
509 tively. In spatial discretization, N_x corresponds to the number of tokens, i.e., the number of
510 channels in CvT. Equipped with data pre- and post-processing as discussed in Section 4.2, the
511 discretized new operator splitting scheme recovers a one-head one-stage CvT.

512 For a K -stage CvT, the number of tokens changes from stage to stage. This property
 513 can be realized by incorporating the new operator-splitting scheme with a hybrid splitting
 514 scheme [42]. Specifically, we need to define K convolution kernels and a bias for convolutional
 515 token embedding. Based on the multigrid idea, these kernels are discretized at different scales,
 516 corresponding to various numbers of tokens.

517 **7. Conclusion.** In this paper, we have introduced a novel operator-theoretic framework
 518 that interprets the Transformer architecture as a discretized operator-splitting scheme for
 519 a continuous integro-differential equation. Within this formulation, core components of the
 520 Transformer—such as self-attention, layer normalization, and feedforward networks—are rig-
 521 orously derived as substeps in a structured variational process and operator-splitting tech-
 522 niques. This approach establishes a unified mathematical foundation that applies broadly to
 523 Transformer-based models, including the original Transformer [46], Vision Transformer (ViT)
 524 [8], and Convolutional vision Transformer (CvT) [50]. By bridging deep learning architectures
 525 with continuous mathematical modeling, our framework advances theoretical understanding
 526 and offers a principled pathway for the design, analysis, and control of future neural architec-
 527 tures. We believe this perspective will stimulate further developments in both the theoretical
 528 analysis and practical refinement of attention-based models.

529 **Appendix.**

530 **Appendix A. Explanation of the variables for Transformers.** The input to the Trans-
 531 former model is usually a sequence. In the context of language processing, the input is usually
 532 sentences. Transformers do not process raw text directly. Instead, they break sentences into
 533 tokens. A token can be a word (e.g., "hello"), a subword (e.g., "multi" + "grid" = "multigrid"),
 534 or even a single character. Then, each different token is transformed into unique embedding
 535 vectors through an embedding layer. The way text is tokenized affects how efficiently a model
 536 understands language.

For example, we have an input sentence "Life is like a boat", and we treat each word as a token. The embedding layer would transform each token into a 512-dimensional vector. Then, the input sentence is represented by a matrix U of size 5×512 :

$$\underbrace{\text{Life is like a boat}}_{\text{Input sentence}} \rightarrow \underbrace{["\text{Life}", "\text{is}", "\text{like}", "\text{a}", "\text{boat}"]}_{\text{Tokens}} \rightarrow U = \underbrace{\begin{bmatrix} \text{Embedding}(\text{"Life"}) \\ \text{Embedding}(\text{"is"}) \\ \text{Embedding}(\text{"like"}) \\ \text{Embedding}(\text{"a"}) \\ \text{Embedding}(\text{"boat"}) \end{bmatrix}}_{\text{Embedded matrix}}.$$

537 In general, a sentence of N_x tokens with N_y -dimensional embedding would be transformed into
 538 a matrix U of size $N_x \times N_y$. This matrix U can be viewed as the discretization of a continuous
 539 function $u(\mathbf{x}, \mathbf{y}, t) : \Omega_x \times \Omega_y \rightarrow \mathbb{R} \times [0, T)$, where Ω_x and Ω_y are continuous domains, on
 540 a structured grid of size $N_x \times N_y$. Here, the variable \mathbf{x} indicates the position of tokens, \mathbf{y}
 541 indicates the entry of token vectors, and t indicates the time variable. The embedded matrix
 542 U is fed into the Transformer model and updated by attention layers.

543 **Appendix B. Proof of Theorem 3.1.**

$$\begin{aligned}
563 \quad & = (u - v) + \frac{1}{|\Omega_y|} \int_{\Omega_y} v(\mathbf{x}, \xi) d\xi - \sigma_1 + \left(-1 + \frac{\int_{\Omega_y} v(\mathbf{x}, \xi)(u(\mathbf{x}, \xi) - \sigma_1) d\xi}{|\Omega_y| \sigma_2^2} \right) u \\
564 \quad & - \left(-1 + \frac{\int_{\Omega_y} v(\mathbf{x}, \xi)(u(\mathbf{x}, \xi) - \sigma_1) d\xi}{|\Omega_y| \sigma_2^2} \right) \sigma_1 \\
565 \quad (B.9) \quad & = \frac{\int_{\Omega_y} v(\mathbf{x}, \xi)(u(\mathbf{x}, \xi) - \sigma_1) d\xi}{|\Omega_y| \sigma_2^2} (u - \sigma_1) - v + \frac{1}{|\Omega_y|} \int_{\Omega_y} v(\mathbf{x}, \xi) d\xi.
\end{aligned}$$

566 We have

$$567 \quad (B.10) \quad \frac{\int_{\Omega_y} v(\mathbf{x}, \xi)(u(\mathbf{x}, \xi) - \sigma_1) d\xi}{|\Omega_y| \sigma_2^2} (u - \sigma_1) = v - \frac{1}{|\Omega_y|} \int_{\Omega_y} v(\mathbf{x}, \xi) d\xi.$$

568 Define $\alpha(\mathbf{x}), \beta(\mathbf{x})$ as in Theorem 3.1. We claim that the solution to (B.10) is given as

$$569 \quad (B.11) \quad u = \frac{v - \alpha}{\sqrt{\beta}} \sigma_2 + \sigma_1.$$

570 To verify it, note that (B.11) can be rewritten as

$$571 \quad (B.12) \quad \frac{u - \sigma_1}{\sigma_2} = \frac{v - \alpha}{\sqrt{\beta}}.$$

572 Utilizing (B.12), we deduce

$$\begin{aligned}
573 \quad & \frac{\int_{\Omega_y} v(\mathbf{x}, \xi)(u(\mathbf{x}, \xi) - \sigma_1) d\xi}{|\Omega_y| \sigma_2^2} (u - \sigma_1) \\
574 \quad & = \frac{1}{|\Omega_y|} \left(\int_{\Omega_y} v(\mathbf{x}, \xi) \frac{(u(\mathbf{x}, \xi) - \sigma_1)}{\sigma_2} d\xi \right) \frac{u - \sigma_1}{\sigma_2} \\
575 \quad & = \frac{1}{|\Omega_y|} \left(\int_{\Omega_y} v(\mathbf{x}, \xi) \frac{(v(\mathbf{x}, \xi) - \sigma_1)}{\sqrt{\beta}} d\xi \right) \frac{v - \sigma_1}{\sqrt{\beta}} \\
576 \quad & = \frac{1}{|\Omega_y| \beta} \left(\int_{\Omega_y} [v^2(\mathbf{x}, \xi) - 2\sigma_1 v(\mathbf{x}, \xi) + \sigma_1^2] d\xi \right) (v - \sigma_1) \\
577 \quad & = \frac{1}{|\Omega_y| \beta} \left(\int_{\Omega_y} (v(\mathbf{x}, \xi) - \sigma_1)^2 d\xi \right) (v - \sigma_1) \\
578 \quad & = \frac{1}{|\Omega_y| \beta} (|\Omega_y| \beta) (v - \sigma_1) \\
579 \quad (B.13) \quad & = v - \frac{1}{|\Omega_y|} \int_{\Omega_y} v(\mathbf{x}, \xi) d\xi,
\end{aligned}$$

580 where the second equality follows from (B.12). The derivation above verifies that (B.11) is the
581 solution to (B.10). ■

582 **Appendix C. An Introduction to Operator Splitting Schemes.**

583 Operator-splitting methods decompose a complicated time-evolution problem into several
584 substeps so that each substep can be solved efficiently. Consider the evolution equation

$$585 \quad (C.1) \quad \begin{cases} u_t + \sum_{k=1}^K A_k(t; u) = 0 \text{ in } \Omega \times (0, T], \\ u(0) = u_0, \end{cases}$$

586 where $A_k(t; u)$'s are operators applied on u . Denote $t^n = n\Delta t$. Given u^n , there are three types
587 of splitting strategies to compute u^{n+1} .

588 Sequential Lie scheme [14, Sec. 2.2] decomposes this computation into K sequential sub-
589 steps. Given u^n , one computes u^{n+1} as follows:

590 For $k = 1, \dots, K$, solve

$$591 \quad (C.2) \quad \begin{cases} u_t + A_k(t; u) = 0 \text{ in } \Omega \times (t^n, t^{n+1}], \\ u(t^n) = u^{n+(k-1)/K}, \end{cases}$$

592 and set $u^{n+k/K} = u(t^{n+1})$.

593 Parallel splitting [32] decomposes the computation into K parallel substeps:

594 For $k = 1, \dots, K$, compute

$$595 \quad (C.3) \quad \begin{cases} v_t + K A_k(t; v) = 0 \text{ in } \Omega \times (t^n, t^{n+1}], \\ v(t^n) = u^n, \end{cases}$$

and set $v_k = u(t^{n+1})$. Then compute

$$u^{n+1} = \frac{1}{K} \sum_{k=1}^K v_k.$$

596 Hybrid splitting [42] is a combination of sequential and parallel splitting, which decomposes
597 the computation into sequential substeps, each of which contains several parallel substeps.

598 Substep (C.2) and (C.3) can be solved by one-step explicit or implicit scheme. For example,
599 (C.2) can be solved by

$$600 \quad \frac{u^{n+k/K} - u^{n+(k-1)/K}}{\Delta t} + A_k(u^{n+k/K}) = 0 \quad \text{or} \quad \frac{u^{n+k/K} - u^{n+(k-1)/K}}{\Delta t} + A_k(u^{n+(k-1)/K}) = 0.$$

601 Our scheme in Section 3.2 is a sequential splitting scheme, in which (3.2) uses explicit scheme,
602 and (3.3) and (3.6) use implicit scheme. Substep (3.4) is semi-implicit, which is slightly
603 different from others. In fact, (3.4) is a composition of two substeps, as it is solved by another
604 sequential splitting (3.15). In (3.15), the first substep is explicit and the second one is implicit.

605 It can be shown that with proper conditions, all three strategies are first-order accurate
606 schemes for the evolution equation.

607

REFERENCES

- 608 [1] J. L. BA, J. R. KIROS, AND G. E. HINTON, *Layer normalization*, arXiv preprint arXiv:1607.06450,
609 (2016).
- 610 [2] M. BENNING, E. CELLEDONI, M. J. EHRHARDT, B. OWREN, AND C.-B. SCHÖNLIEB, *Deep learning*
611 *as optimal control problems: Models and numerical methods*, Journal of Computational Dynamics, 6
612 (2019), pp. 171–198.
- 613 [3] A. BRYUTKIN, J. HUANG, Z. DENG, G. YANG, C.-B. SCHÖNLIEB, AND A. AVILES-RIVERO, *Ham-*
614 *let: graph transformer neural operator for partial differential equations*, in Proceedings of the 41st
615 International Conference on Machine Learning, 2024, pp. 4624–4641.
- 616 [4] M. CHEN, H. JIANG, W. LIAO, AND T. ZHAO, *Nonparametric regression on low-dimensional manifolds*
617 *using deep relu networks: Function approximation and statistical recovery*, Information and Inference:
618 A Journal of the IMA, 11 (2022), pp. 1203–1253.
- 619 [5] R. T. CHEN, Y. RUBANOVA, J. BETTENCOURT, AND D. K. DUVENAUD, *Neural ordinary differential*
620 *equations*, Advances in neural information processing systems, 31 (2018).
- 621 [6] C.-W. CHENG, J. HUANG, Y. ZHANG, G. YANG, C.-B. SCHÖNLIEB, AND A. I. AVILES-RIVERO,
622 *Mamba neural operator: Who wins? transformers vs. state-space models for pdes*, arXiv preprint
623 arXiv:2410.02113, (2024).
- 624 [7] Y. CUI, Y. XU, R. PENG, AND D. WU, *Layer normalization for tsk fuzzy system optimization in*
625 *regression problems*, IEEE Transactions on Fuzzy Systems, 31 (2022), pp. 254–264.
- 626 [8] A. DOSOVITSKIY, L. BEYER, A. KOLESNIKOV, D. WEISSENBORN, X. ZHAI, T. UNTERTHINER, M. DE-
627 HGHANI, M. MINDERER, G. HEIGOLD, S. GELLY, ET AL., *An image is worth 16x16 words: Trans-*
628 *formers for image recognition at scale*, in International Conference on Learning Representations, 2021.
- 629 [9] S. DUTTA, T. GAUTAM, S. CHAKRABARTI, AND T. CHAKRABORTY, *Redesigning the transformer architec-*
630 *ture with insights from multi-particle dynamical systems*, Advances in Neural Information Processing
631 Systems, 34 (2021), pp. 5531–5544.
- 632 [10] T. FURUYA, M. V. DE HOOP, AND G. PEYRÉ, *Transformers are universal in-context learners*, arXiv
633 preprint arXiv:2408.01367, (2024).
- 634 [11] B. GESHKOVSKI, C. LETROUIT, Y. POLYANSKIY, AND P. RIGOLLET, *A mathematical perspective on*
635 *transformers*, arXiv preprint arXiv:2312.10794, (2023).
- 636 [12] R. GLOWINSKI AND P. LE TALLEC, *Augmented Lagrangian and operator-splitting methods in nonlinear*
637 *mechanics*, vol. 9, Society for Industrial Mathematics, 1989.
- 638 [13] R. GLOWINSKI, S. J. OSHER, AND W. YIN, *Splitting methods in communication, imaging, science, and*
639 *engineering*, Springer, 2017.
- 640 [14] R. GLOWINSKI, T.-W. PAN, AND X.-C. TAI, *Some facts about operator-splitting and alternating direction*
641 *methods*, in Splitting Methods in Communication, Imaging, Science, and Engineering, Springer, 2016,
642 pp. 19–94.
- 643 [15] A. GRAVES, A.-R. MOHAMED, AND G. HINTON, *Speech recognition with deep recurrent neural networks*,
644 in 2013 IEEE international conference on acoustics, speech and signal processing, Ieee, 2013, pp. 6645–
645 6649.
- 646 [16] K. GREGOR AND Y. LECUN, *Learning fast approximations of sparse coding*, in Proceedings of the 27th
647 international conference on international conference on machine learning, 2010, pp. 399–406.
- 648 [17] E. HABER, L. RUTHOTTO, E. HOLTHAM, AND S.-H. JUN, *Learning across scales—multiscale methods for*
649 *convolution neural networks*, in Proceedings of the AAAI conference on artificial intelligence, vol. 32,
650 2018.
- 651 [18] P. HAGEMANN, J. HERTRICH, AND G. STEIDL, *Stochastic normalizing flows for inverse problems: A*
652 *markov chains viewpoint*, SIAM/ASA Journal on Uncertainty Quantification, 10 (2022), pp. 1162–
653 1190.
- 654 [19] A. HAVRILLA AND W. LIAO, *Understanding scaling laws with statistical and approximation theory for*
655 *transformer neural networks on intrinsically low-dimensional data*, Advances in Neural Information
656 Processing Systems, 37 (2024), pp. 42162–42210.
- 657 [20] J. HE AND J. XU, *Mgnet: A unified framework of multigrid and convolutional neural network*, Science
658 china mathematics, 62 (2019), pp. 1331–1354.
- 659 [21] S. JELASSI, M. SANDER, AND Y. LI, *Vision transformers provably learn spatial structure*, Advances in
660 Neural Information Processing Systems, 35 (2022), pp. 37822–37836.
- 661 [22] F. JIANG, Y. JIANG, H. ZHI, Y. DONG, H. LI, S. MA, Y. WANG, Q. DONG, H. SHEN, AND Y. WANG,

- 662 *Artificial intelligence in healthcare: past, present and future*, Stroke and vascular neurology, 2 (2017).
- 663 [23] A. KRIZHEVSKY, I. SUTSKEVER, AND G. E. HINTON, *Imagenet classification with deep convolutional*
664 *neural networks*, Advances in neural information processing systems, 25 (2012).
- 665 [24] Z. LAI, L.-H. LIM, AND Y. LIU, *Attention is a smoothed cubic spline*, arXiv preprint arXiv:2408.09624,
666 (2024).
- 667 [25] Y. LAN, Z. LI, J. SUN, AND Y. XIANG, *Dosnet as a non-black-box pde solver: When deep learning meets*
668 *operator splitting*, Journal of Computational Physics, 491 (2023), p. 112343.
- 669 [26] Z. LI, N. KOVACHKI, K. AZIZZADENESHELI, B. LIU, K. BHATTACHARYA, A. STUART, AND
670 A. ANANDKUMAR, *Fourier neural operator for parametric partial differential equations*, arXiv preprint
671 arXiv:2010.08895, (2020).
- 672 [27] H. LIU, J. LIU, R. H. CHAN, AND X.-C. TAI, *Double-well net for image segmentation*, Multiscale
673 Modeling & Simulation, 22 (2024), pp. 1449–1477.
- 674 [28] Y. LIU, Z. ZHANG, AND H. SCHAEFFER, *Prose: Predicting multiple operators and symbolic expressions*
675 *using multimodal transformers*, Neural Networks, 180 (2024), p. 106707.
- 676 [29] Z. LONG, Y. LU, X. MA, AND B. DONG, *Pde-net: Learning pdes from data*, in International conference
677 on machine learning, PMLR, 2018, pp. 3208–3216.
- 678 [30] J. LU, Z. SHEN, H. YANG, AND S. ZHANG, *Deep network approximation for smooth functions*, SIAM
679 Journal on Mathematical Analysis, 53 (2021), pp. 5465–5506.
- 680 [31] L. LU, P. JIN, G. PANG, Z. ZHANG, AND G. E. KARNIADAKIS, *Learning nonlinear operators via*
681 *deepnet based on the universal approximation theorem of operators*, Nature machine intelligence, 3
682 (2021), pp. 218–229.
- 683 [32] T. LU, P. NEITTAANMAKI, AND X.-C. TAI, *A parallel splitting-up method for partial differential equa-*
684 *tions and its applications to navier-stokes equations*, ESAIM: Mathematical Modelling and Numerical
685 Analysis, 26 (1992), pp. 673–708.
- 686 [33] Y. LU, Z. LI, D. HE, Z. SUN, B. DONG, T. QIN, L. WANG, AND T.-Y. LIU, *Understanding and*
687 *improving transformer from a multi-particle dynamic system point of view.*, in ICLR 2020 Workshop
688 on Integration of Deep Neural Models and Differential Equations.
- 689 [34] S. MARTIN, A. GAGNEUX, P. HAGEMANN, AND G. STEIDL, *Pnp-flow: Plug-and-play image restoration*
690 *with flow matching*, arXiv preprint arXiv:2410.02423, (2024).
- 691 [35] R. MIOTTO, F. WANG, S. WANG, X. JIANG, AND J. T. DUDLEY, *Deep learning for healthcare: review,*
692 *opportunities and challenges*, Briefings in bioinformatics, 19 (2018), pp. 1236–1246.
- 693 [36] M. RAISSI, P. PERDIKARIS, AND G. E. KARNIADAKIS, *Physics-informed neural networks: A deep learn-*
694 *ing framework for solving forward and inverse problems involving nonlinear partial differential equa-*
695 *tions*, Journal of Computational physics, 378 (2019), pp. 686–707.
- 696 [37] O. RONNEBERGER, P. FISCHER, AND T. BROX, *U-net: Convolutional networks for biomedical image*
697 *segmentation*, in Medical image computing and computer-assisted intervention–MICCAI 2015: 18th
698 international conference, Munich, Germany, October 5-9, 2015, proceedings, part III 18, Springer,
699 2015, pp. 234–241.
- 700 [38] D. RUIZ-BALET AND E. ZUAZUA, *Neural ode control for classification, approximation, and transport*,
701 SIAM Review, 65 (2023), pp. 735–773.
- 702 [39] L. RUTHOTTO AND E. HABER, *Deep neural networks motivated by partial differential equations*, Journal
703 of Mathematical Imaging and Vision, 62 (2020), pp. 352–364.
- 704 [40] Z. SHEN, A. HAVRILLA, R. LAI, A. CLONINGER, AND W. LIAO, *Transformers for learning on noisy and*
705 *task-level manifolds: Approximation and generalization insights*, arXiv preprint arXiv:2505.03205,
706 (2025).
- 707 [41] R. STRUDEL, R. GARCIA, I. LAPTEV, AND C. SCHMID, *Segmenter: Transformer for semantic segmenta-*
708 *tion*, in Proceedings of the IEEE/CVF international conference on computer vision, 2021, pp. 7262–
709 7272.
- 710 [42] X.-C. TAI, H. LIU, AND R. CHAN, *Pottsmgnet: A mathematical explanation of encoder-decoder based*
711 *neural networks*, SIAM Journal on Imaging Sciences, 17 (2024), pp. 540–594.
- 712 [43] X.-C. TAI, H. LIU, R. H. CHAN, AND L. LI, *A mathematical explanation of unet*, Mathematical
713 Foundations of Computing, (2024).
- 714 [44] S. TAKAKURA AND T. SUZUKI, *Approximation and estimation ability of transformers for sequence-to-*
715 *sequence functions with infinite dimensional input*, in International Conference on Machine Learning,

- 716 PMLR, 2023, pp. 33416–33447.
- 717 [45] R. E. TURNER, *An introduction to transformers*, arXiv preprint arXiv:2304.10557, (2023).
- 718 [46] A. VASWANI, N. SHAZEER, N. PARMAR, J. USZKOREIT, L. JONES, A. N. GOMEZ, Ł. KAISER, AND
- 719 I. POLOSUKHIN, *Attention is all you need*, Advances in neural information processing systems, 30
- 720 (2017).
- 721 [47] T. WANG, Z. DOU, C. BAO, AND Z. SHI, *Diffusion mechanism in residual neural network: theory and*
- 722 *applications*, IEEE Transactions on Pattern Analysis and Machine Intelligence, 46 (2023), pp. 667–
- 723 680.
- 724 [48] E. WEINAN, *A proposal on machine learning via dynamical systems*, Communications in Mathematics
- 725 and Statistics, 5 (2017), pp. 1–11.
- 726 [49] E. WEINAN, C. MA, AND L. WU, *Machine learning from a continuous viewpoint, i*, Science China
- 727 Mathematics, 63 (2020), pp. 2233–2266.
- 728 [50] H. WU, B. XIAO, N. CODELLA, M. LIU, X. DAI, L. YUAN, AND L. ZHANG, *CvT: Introducing convolu-*
- 729 *tions to vision transformers*, in Proceedings of the IEEE/CVF international conference on computer
- 730 vision, 2021, pp. 22–31.
- 731 [51] J. XU, X. SUN, Z. ZHANG, G. ZHAO, AND J. LIN, *Understanding and improving layer normalization*,
- 732 Advances in neural information processing systems, 32 (2019).
- 733 [52] Y. YANG, J. SUN, H. LI, AND Z. XU, *Admm-csnet: A deep learning approach for image compressive*
- 734 *sensing*, IEEE transactions on pattern analysis and machine intelligence, 42 (2018), pp. 521–538.
- 735 [53] D. YAROTSKY, *Error bounds for approximations with deep relu networks*, Neural networks, 94 (2017),
- 736 pp. 103–114.
- 737 [54] Z. YE, X. HUANG, L. CHEN, H. LIU, Z. WANG, AND B. DONG, *Pdeformer: Towards a foundation*
- 738 *model for one-dimensional partial differential equations*, arXiv preprint arXiv:2402.12652, (2024).
- 739 [55] T. YOUNG, D. HAZARIKA, S. PORIA, AND E. CAMBRIA, *Recent trends in deep learning based natural*
- 740 *language processing*, iee Computational intelligence magazine, 13 (2018), pp. 55–75.
- 741 [56] S. YUN, M. JEONG, R. KIM, J. KANG, AND H. J. KIM, *Graph transformer networks*, Advances in
- 742 neural information processing systems, 32 (2019).
- 743 [57] B. ZHANG AND R. SENNRICH, *A lightweight recurrent network for sequence modeling*, in Proceedings of
- 744 the 57th Annual Meeting of the Association for Computational Linguistics, 2019, pp. 1538–1548.
- 745 [58] D.-X. ZHOU, *Universality of deep convolutional neural networks*, Applied and computational harmonic
- 746 analysis, 48 (2020), pp. 787–794.
- 747 [59] J. ZHU, X. CHEN, K. HE, Y. LECUN, AND Z. LIU, *Transformers without normalization*, arXiv preprint
- 748 arXiv:2503.10622, (2025).
- 749 [60] A. ZIAEE AND E. ÇANO, *Batch layer normalization a new normalization layer for cnns and rnns*, in
- 750 Proceedings of the 6th International Conference on Advances in Artificial Intelligence, 2022, pp. 40–
- 751 49.