

---

# Week 1 Reading Summaries: Chapter 1 and 4

---

**Amogh Patankar**

University of California, San Diego  
3869 Miramar Street Box #3037, La Jolla, CA- 92092  
apatankar@ucsd.edu

## Abstract

The concept of Machine Learning Systems is a special concept that is the bridge between theoretical foundations of AI/ML, and the hands-on, practical engineering required to implement those theories. These notes serve as an introductory piece, and a deep dive into DNN Architectures.

## 1 Chapter 1: Introduction

### 1.1 AI is Everywhere

AI is a superhuman force that is transforming how we live, from infrastructure and healthcare, to more scientific research. The AI Revolution is advancing at an unprecedented pace, and the potential is massive- potentially life changing at a global scale. This rapid evolution is an engineering challenge, i.e. developing systems that can achieve superhuman capabilities in many domains.

### 1.2 Understanding AI/ML

Intelligent systems are the key to AI, and the replication of human behavior, learning, adapting, and reasoning, is the fundamental problem at hand. ML is used to build these systems, and the relationship with AI is similar to the connection between theory and practice. This progression drives the creation of systems that are similar to intelligent behavior.

### 1.3 Evolution of AI

AI has changed in the past 7 decades, starting with the perception in 1957, all the way to GPT4 in 2023. The point being, the early AI systems don't even compare to the model ML systems that do a wide variety of applications and tasks at the present day.

#### 1.3.1

AI was coined at the Dartmouth Conference in 1956, based on the idea that intelligence is equal to symbol manipulation, thanks to the STUDENT system. STUDENT and similar AI could only process input that they stored. The limitation revealed a deep problem, that AI at the time could only match patterns.

#### 1.3.2

MYCIN was developed at Stanford, to diagnose blood infections in the 70s. Again, we saw challenges- that domain knowledge and converting it to rules was difficult, and doctors for example, could process better than the model. The challenges of knowledge and handling uncertainty was something that is still a concern in modern ML.

### 30 1.3.3

31 Statistical learning was built off the premise of Moore's Law, that algorithms like SVMs and NNs  
32 could learn patterns because of a larger amount of data and compute. Statistical learning changed the  
33 approach towards building AI due to: quality/quantity of data, evolutionary, and precision and recall  
34 tension.

### 35 1.3.4

36 Shallow learning (a precursor to DL) was due to powerful algorithms like decision trees and KNNs  
37 process data well, but was limited to small model sizes. Yet, they were relevant for nearly a decade.

### 38 1.3.5

39 Deep learning was the crux of modern ML, relying on larger layers of neurons, with each layer  
40 transforming data to different representations. Breakthroughs like AlexNet and others led to more and  
41 more powerful models being developed, off the premise of complex, deep networks. This foundation  
42 led to unprecedented scaling, with models 1000x the size of AlexNet. The progression from symbolic  
43 reasoning to DL allowed scaling, and for applications to become more complex.

## 44 1.4 Rise of MLSysEng

45 The shift in the past 10 years built on the idea that engineering was required to build strong ML, i.e.  
46 ML systems. The idea of advancing HW, SW, EE, and ML together created MLSysEng, i.e. the  
47 interplay between algorithms and engineering to advance ML.

## 48 1.5 Definition of MLSys

49 A ML System is defined as an integrated system constituting of data, learning algorithms, and  
50 computing infra. Each aspect enhances/dictates the other two, and the interdependency of these  
51 aspects indicates no independence.

## 52 1.6 MLSys Lifecycle

53 The lifecycle of ML Systems consists of data collection and preprocessing, model training and  
54 evaluation, deployment and, finally, monitoring. These interconnected stages allow for smooth and  
55 reliable improvements in model performance, and systems. Each piece limits the improvement of the  
56 system, and problems can compound.

## 57 1.7 Spectrum of MLSys

58 ML Systems range from tinyML systems running on microcontrollers, and scale all the way up to  
59 cloud-based LLM Systems in data centers. Regardless, there is a broad spectrum of ML systems,  
60 with various capabilities depending on size, compute requirements, resources, and limitations of  
61 applications.

## 62 1.8 MLSys Implications on ML Lifecycle

63 A bottle neck of the ML lifecycle is ML Systems- hardware often time determines architectural  
64 decisions, and resource management varies drastically. Resources considerations influence both  
65 model and system architecture, and complexity increases with system distribution. The tradeoffs of  
66 compute, architecture, complexity, and systems requires a delicate balance.

### 67 1.8.1 Emerging Trends

68 There is innovation in application, as well as system architecture. App-level innovation can handle  
69 system changes, allowing advancements. System architecture changes require resource efficiency to  
70 advance, and hardware such as AI accelerators to advance as well. Both advancements allow systems  
71 to grow sustainably.

## 72 **1.9 Real-World Applications**

73 There are diverse ML system applications.

### 74 **1.9.1 FarmBeats**

75 ML in non-tech industries requires lots of specifications to succeed. For example, in FarmBeats, the  
76 data ecosystem is diverse, the ML models are tailored to agricultural applications, and the compute  
77 allows models to run on [small] IoT and edge devices. All in all, FarmBeats showed great deployment  
78 of ML systems in real-world environments.

### 79 **1.9.2 AlphaFold**

80 In the case of AlphaFold, the same aspects (data, algorithm, and infra/compute) were a point of  
81 success. A complex data pipeline with preprocessing, novel NN architectures, and exceptional  
82 parallel computing resources allowed for true innovation.

### 83 **1.9.3 Autonomous Vehicles**

84 For Waymo, the data was of a different manner- namely, data from LiDAR sensors, generating  
85 1 Tb/hour. The structured and unstructured nature of this data was the unique aspect, along with  
86 its algorithm and compute, that allowed Waymo to be effective, and to potentially revolutionize  
87 transportation and urban planning.

## 88 **1.10 Challenges + Considerations**

89 There are many challenges in ML Systems that are important.

### 90 **1.10.1 Data Challenges**

91 Data management in ML Systems poses a unique challenge- as high quality and large quantity data  
92 is hard to obtain, leading to techniques like synthetic data generation. ML Systems require large  
93 amounts of data, and this poses storage/processing challenges as well, alongside data drift.

### 94 **1.10.2 Model Challenges**

95 ML models, which can be as complex as hundreds of billions of parameters, create problems such as  
96 situations with limited resources. The training process has trade-offs which are a challenge, especially  
97 in real world conditions.

### 98 **1.10.3 System Challenges**

99 System challenges include training and serving systems being up to par, meaning that monitoring and  
100 updating must also be smoothly carried out.

### 101 **1.10.4 Ethical + Social Considerations**

102 Fairness and transparency in ML Systems is a long-standing problem that researchers have only  
103 started to dive into and make advances in. Privacy, especially in the high tech modern world, is a  
104 major concern, especially with data and systems themselves.

## 105 **1.11 Future Directions**

106 The future of ML Systems is bright, and trends such as democratization of AI, efficient ML, and  
107 autonomous ML systems, are the future of the field. As such, the limitations posed by AI and ML  
108 systems are difficult to tackle, leading to more research in the future.

## 109 2 Chapter 4: DNN Architectures

### 110 2.1 Overview

111 Deep learning architectures are how models that we use in deep learning are organized- the depth,  
112 breadth, complexity, and structure. Neural networks have evolved from simple pattern recognition,  
113 all the way to generating text, audio and video, due to more complex architectures. Architectures are  
114 usually discussed w.r.t. their algorithmic basics, but each pattern maps differently to HW resources.  
115 Specifically, memory access patterns, computation characteristics, data movement, and resource  
116 allocation.

### 117 2.2 Multi-Layer Perceptrons

118 MLPs are the "easiest" deep neural net architecture, consisting of mathematical computations (matrix  
119 multiplications), and systematically progressing through layers.

#### 120 2.2.1 Pattern Processing Needs

121 DL Systems run into problems where input is directly influential on output, and the solution is dense  
122 pattern processing. This allows for output to be somewhat independent of input combinations, and  
123 feature importance is introduced, allowing systems to determine the most important connections.  
124 Moreover, adaptive representations lets the network(s) reshape representations based on input data.

#### 125 2.2.2 Algorithmic Structure

MLPs have fully connected layers, and this translates to MatMul operations. Essentially, the formula  
looks like this:  $h^l = f(W^l n^{l-1} + b^l)$

126 This creates patterns to be handled efficiently by computer systems.

#### 127 2.2.3 Computational Mapping

128 There are a few layers to computational mapping. The first computation is the simple mathematical  
129 equation in the previous section, i.e. the matmul of X and W, with a bias, fed into an activation.  
130 The second computation involves batch processing, computing neuron outputs, and passing them  
131 through an activation function. This translation from mathematics to computation shows the desnity  
132 of matmul operations.

#### 133 2.2.4 System Implications

134 Three fundamental dimensions exist for computation patterns- memory requirements, computational  
135 needs, and data movement, which enables analysis of system design decisions. With respect to  
136 memory requirements, the weights, inputs, and results impact memory requirements. Optimization  
137 opportunities present themselves through data organization/reuse. For computation, MAC (multiply-  
138 accumulate) operations are nested; this lends itself to optimization strategies, such as BLAS. Finally,  
139 data movement allows for data staging and transfer optimization.

### 140 2.3 CNN: Spatial Pattern Processing

141 CNNs are the best models for encapsulating spatial relationships, namely, pixels. As such, CNNs  
142 have different needs and approaches than MLPs.

#### 143 2.3.1 Pattern Processing Needs

144 Spatial processing is specifically useful for data depending on relative positions, such as an edge in  
145 an image. So, we require the patterns to retain their structure/"shape", i.e. pattern localization and  
146 pattern recognition, independent of position. CNNs use local connections to achieve this, known as  
147 convolution (set of weights applied).

### 148 2.3.2 Algorithmic Structure

$$\mathbf{H}_{i,j,k}^{(l)} = f\left(\sum_{di} \sum_{dj} \sum_c \mathbf{W}_{di,dj,c,k}^{(l)} \mathbf{H}_{i+di,j+dj,c}^{(l-1)} + \mathbf{b}_k^{(l)}\right)$$

Figure 1: With respect to CNNs, the core operation looks like this.

149 This operation process local regions, reuses the same weights, and maintains spatial structure,  
150 differing from the MLP architecture. This convolution operation moves throughout the image as a  
151 filter.

### 152 2.3.3 Computational Mapping

153 The first map is a convolution operation, i.e using kernel, input, bias, and an activation. The second  
154 map applies the same weight filters to local regions. The result is a  $m \times m$  filter, with  $n$  output  
155 channels, with significantly less MAC ops than MLP.

### 156 2.3.4 System Implications

157 Here we look at- memory requirements, computational needs, and data movement. CNNs use reusable  
158 filters, unlike MLPs, and hence require much less memory, as we reuse weights and optimize in that  
159 manner. With regards to computational needs, the load remains large due to spatial repetition despite  
160 fewer ops than MLPs. Finally, the window of convolutions and reused weights allows for streaming  
161 input.

## 162 2.4 RNN: Sequential Pattern Processing

163 While MLPs are tailored to handle arbitrary relationships, and CNNs handle spatial, RNNs are meant  
164 for sequential data [patterns].

### 165 2.4.1 Pattern Processing Needs

166 For RNNs, the primary use case is NLP, meaning that a word's meaning is dependent on the previous  
167 sequence. That is, the challenge is maintaining context, and RNNs have variable-length sequences.

### 168 2.4.2 Algorithmic Structure

169 RNNs use a different structure, by using "recurrent connections", allowing RNNs to maintain a state  
170 that gets updated. Sequential processing uses a high-dimensional vector for each word, and uses an  
171 equation like the one that follows:

$$\mathbf{h}_t = f(\mathbf{W}_{hh}\mathbf{h}_{t-1} + \mathbf{W}_{xh}\mathbf{x}_t + \mathbf{b}_h)$$

Figure 2

### 172 2.4.3 Computational Mapping

173 Computational mapping for RNNs is a single time step, with input, hidden state, and two weight  
174 matrices, passed through an activation function. Nested loops then process each sequence (batch  
175 parallelism), previous hidden states, new information, and biases + activations.

### 176 2.4.4 System Implications

177 For RNNs, the memory requirements, computation needs, and data movement differ than that of  
178 a MLP or a CNN. We store two sets of weights, and a hidden state, meaning memory usage is  
179 heightened. By extension, computational needs cannot parallelize across steps due to sequential  
180 nature, and similarly, weights can be reused for data flow.

## 181 2.5 Attention Mechanism

182 The attention mechanism was specifically built for language understanding, i.e. to capture relation-  
183 ships between words rather than just position.

### 184 2.5.1 Pattern Processing Needs

185 The attention architecture can dynamically process relationships between two words, and systems  
186 must be able to compute relationships, weigh them, and use weights to combine information.

### 187 2.5.2 Basic Attention Mechanism

- 188 1. **Algorithmic Structure:** The real power of the attention mechanism is computed by using  
189 queries, keys, and values, representing learned projections of input. The calculation is  
190  $\text{Attention}(Q, K, V) = \text{softmax}((QK^T * d_k^{-1/2}) * V)$ . The Q, K, V projection get calculated,  
191 an N x N attention matrix is formed, and we combine value vectors, using attention weights,  
192 generating an output.
- 193 2. **Computational Mapping:** The mapping is the mathematical abstraction described above,  
194 with nested loops for sequence processing, attention compute, key comparison, and attention  
195 score.
- 196 3. **System Implications:** For the attention architecture, there must be storage for attention  
197 weights, KQV projections, and feature representations. Similarly, the system performs an  
198 enormous amount of MAC operations, and by extension, dynamically computed weights  
199 must move frequently, requiring high data movement.

### 200 2.5.3 Transformers and Self-Attention

201 Transformers extend the idea of attention by applying it for a single sequence, resulting in one element  
202 "attending" to all others.

- 203 1. **Algorithmic Structure:** The innovative aspect of transformers is due to self-attention layers,  
204 in which QKV are determined from the same input, meaning the weights can be calculated  
205 while encoding positions. This allows it to capture long-range dependencies; transformers  
206 generally use multi-head attention, with multiple sets of QKV projections.
- 207 2. **Computational Mapping:** The computational mapping of attention is increased drastically,  
208 due to the introduction of self-attention, and multi-head attention simultaneously.
- 209 3. **System Implications:** Self-attention allows for parallel processing for all positions; the  
210 attention score calculation has quadratic complexity w.r.t sequence length, which is a massive  
211 bottleneck. Moreover, the multi-head attention is a parallel run of self-attention ops, which  
212 linearly increases the power. And finally, the computations themselves, and intermediate  
213 results are all very memory heavy and dominant.

## 214 2.6 Architectural Building Blocks

### 215 2.6.1 Perceptron to MLP

216 The transition from perceptron to MLP was a new idea that is in every new modern architecture, such  
217 as the MLP feed forward style networks in transformers. The idea of data transformation through  
218 non-linear layers and backprop, is a building block of modern ML.

### 219 2.6.2 Dense to Spatial Processing

220 CNNs revolutionized data processing, with the ideas of parameter sharing, and skip connections.  
221 These ideas allowed for more efficient networks, but also direct paths for gradient flow, and of course,  
222 batch normalization.

223 **2.6.3 Evolution of Sequence Processing**

224 LSTMs and GRUs allowed gating, which eventually lead to a pattern in architecture design, with  
 225 the idea that networks could process variable-length inputs by reusing weights, and eventually led to  
 226 attention.

227 **2.6.4 Modern Architectures: Synthesis and Innovation**

228 Modern architectures are built from an evolution of NNs, from MLPs to CNNs to RNNs and finally  
 229 using attention in the modern era. The ideas of self-attention (allowing for dynamic processing) and  
 230 FF with skip and normalization, has changed modern language modeling.

231 And of course, with change, comes various utilization ideas, such as this:

Primitive Type	MLP	CNN	RNN	Transformer
Computational	Matrix Multiplication	Convolution (Matrix Mult.)	Matrix Mult. + State Update	Matrix Mult. + Attention
Memory Access	Sequential	Strided	Sequential + Random	Random (Attention)
Data Movement	Broadcast	Sliding Window	Sequential	Broadcast + Gather

Figure 3

233 **2.7 System-Level Building**

234 **2.7.1 Memory Access Primitives**

235 With the newest architectures, we must also look at various primitives for memory access. Sequential  
 236 access allows for memory like DRAM to read quickly, while strided access for CNNs requires  
 237 software frameworks to change this data to sequential access. Random access posing the largest  
 238 threat for efficiency is the reason why LMSys is such a large domain to work and research in.

239 **2.7.2 Data Movement Primitives**

240 When we talk about data movement, we look for on-chip vs off-chip, as one requires much more  
 241 energy than the other. Broadcast, scatter, gather, and reduce are key to distribute and collect data  
 242 across computational units. Broadcasting- sends data to many destinations, scatter- distributes  
 243 different data to different locations, gather- collects from multiple sources, and reduce combines  
 244 multiple values to one result.

245 **2.7.3 System Design Impact**

246 The effect of all these primitives leads us to explore different designs for systems in ML. Operations  
 247 like matmul, convolution, and MAC lead to changing of processing units. Balancing the tradeoff  
 248 needs us to evaluate the workload and deployment, to choose the right memory and data primitives to  
 249 make good decisions.