

DR.-ING. MARK SCHUTERA

AGENTIC AI AND AGENT SYSTEMS

UNFINISHED LECTURE NOTES

Copyright © 2025 Dr.-Ing. Mark Schutera

PUBLISHED BY UNFINISHED LECTURE NOTES

Licensed under the Creative Commons Attribution-NonCommercial 4.0 International License (“CC BY-NC 4.0”). You may not use this file for commercial purposes. You must obtain explicit permission from the author for uses beyond those permitted by this license. To view a copy of this license, visit <https://creativecommons.org/licenses/by-nc/4.0/>. Unless required by applicable law or agreed to in writing, distributed material is provided on an “AS IS” BASIS, WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied. See the license for details.

First printing, July 2025

THE HOTTEST NEW PROGRAMMING LANGUAGE IS ENGLISH.

ANDREJ KARPATHY

Contents

<i>Foundations of Agentic AI</i>	9
----------------------------------	---

Introduction

These unfinished lecture notes analyze agent systems founded on Large Language Models, examining their theoretical foundations, architectural design, and operational mechanisms. We investigate how LLMS function as cognitive substrates for autonomous agents and their integration of reasoning, action and observation capabilities through structured interaction interfaces.

THESE NOTES ARE A LIVING THING, subject to continuous updates and improvements—and have only recently seen the light of day. They are not yet complete and are rather intended to showcase the look and feel of my lecture notes, but I hope you will come to find them useful nonetheless.

Foundations of Agentic AI

LANGUAGE MODEL EVOLUTION proceeded through discrete developmental phases culminating in contemporary agentic systems. Pre-2017 methodologies employed statistical models and distributed representations—Word2Vec¹ and GloVe²—alongside primitive conversational systems like ELIZA³.

The Transformer architecture⁴ introduced in 2017 revolutionized NLP through self-attention mechanisms, enabling efficient parallel processing and long-range dependency modeling.

The Transformer’s self-attention mechanism⁵ facilitated Large Language Model development (2018-2021), characterized by parametric scale and sophisticated training paradigms: pre-training, few-shot learning, and reinforcement learning from human feedback.

THE 2022 WATERSHED marked by ChatGPT’s deployment transformed LLMS from research artifacts to mainstream applications. GPT-3’s⁶ emergent reasoning capabilities established language models as inference engines.

Agent Formalization and Autonomy

AN AGENT constitutes a computational system leveraging an AI model to reason, plan, and interact with environments through external tool execution to achieve specified objectives.

The distinguishing characteristic of *agency* represents the capacity for autonomous action based on internal reasoning processes. This formalization encompasses both reactive and deliberative behavioral patterns.

THIS TAXONOMIC PROGRESSION from passive processing to autonomous workflow orchestration represents escalating agency levels. LEVEL 4 autonomy emerges when agents initiate and control execution cycles, including meta-agent spawning and coordination.

¹ Word2Vec pioneered distributed word representations, enabling semantic similarity computation through vector arithmetic operations.

T. Mikolov, K. Chen, G. Corrado, and J. Dean. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*, 2013. URL <https://arxiv.org/abs/1301.3781>

² J. Pennington, R. Socher, and C. D. Manning. Glove: Global vectors for word representation. *EMNLP*, 2014. URL <https://nlp.stanford.edu/pubs/glove.pdf>

³ J. Weizenbaum. Eliza—a computer program for the study of natural language communication between man and machine. *Communications of the ACM*, 9(1):36–45, 1966. DOI: 10.1145/365153.365168. URL <https://doi.org/10.1145/365153.365168>

⁴ A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Łukasz Kaiser, and I. Polosukhin. Attention is all you need. *Advances in Neural Information Processing Systems*, 30:5998–6008, 2017. URL <https://papers.nips.cc/paper/7181-attention-is-all-you-need.pdf>

⁵ Attention enables models to weigh the relevance of different input elements when processing each token, fundamentally improving context understanding.

⁶ T. B. Brown, B. Mann, N. Ryder, et al. Language models are few-shot learners. *arXiv preprint arXiv:2005.14165*, 2020. URL <https://arxiv.org/abs/2005.14165>

Full autonomy manifests when agents can recursively instantiate other agents, creating hierarchical control structures that mirror biological neural organization.

Autonomy	Description	Implementation Pattern
LEVEL 0	Passive output processing	<code>process_llm_output(llm_response)</code>
LEVEL 1	Conditional control flow	<code>if llm_decision(): path_a() else: path_b()</code>
LEVEL 2	Dynamic function invocation	<code>run_function(llm_chosen_tool, llm_args)</code>
LEVEL 3	Iterative process control	<code>while llm_should_continue(): execute_step()</code>
LEVEL 4	Meta-agent orchestration	<code>if llm_trigger(): spawn_agent()</code>

Table 1: Hierarchical Agent Autonomy Levels

Language Models as Cognitive Substrates

LARGE LANGUAGE MODELS constitute the foundational cognitive architecture for modern agent systems through their capacity for pattern recognition, contextual reasoning, and linguistic generation. Training on vast textual corpora enables LLMs to internalize linguistic patterns, semantic relationships, and domain-specific knowledge structures⁷. This emergent capability transforms static language models into dynamic reasoning engines.

TOKENIZATION represents the fundamental computational unit whereby LLMs decompose linguistic input into discrete, processable elements—typically words or subword units.

SELF-SUPERVISED LEARNING exploits implicit textual structure, enabling training on unlabeled corpora:

$$\mathcal{D} = \{(x_t, \tilde{y}_t)\}_{t=1}^N \quad \text{where} \quad x_{t+1} = \tilde{y}_t$$

THE ATTENTION MECHANISM computes similarity scores between queries and keys, then applies these weights to corresponding values:

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax} \left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d_k}} \right) \mathbf{V}$$

Attention enables the model to capture **long-range** dependencies within text sequences. The **query** of each individual token determines which other tokens receive higher attention weights. This mechanism supports efficient **batch processing** on GPUs. The output

⁷ T. B. Brown, B. Mann, N. Ryder, et al. Language models are few-shot learners. *arXiv preprint arXiv:2005.14165*, 2020. URL <https://arxiv.org/abs/2005.14165>

Token vocabularies typically exceed 30,000 elements, balancing computational efficiency with semantic granularity.

Figure 1: **Playground**–Tokenization process by merging in Corpora.

Query (Q): Things I am looking for (adjectives).

Key (K): What do I contain? (adjective).

Value (V): What information is stored? (“awesome”).

of the attention layer is a **contextualized representation** of the input tokens, reflecting their relationships in context.

The attention operation computes how well each query matches the keys, then uses these scores to aggregate the corresponding values, producing a context-aware representation for each token.

The normalization factor $\sqrt{d_k}$ prevents dot products from growing excessively large, which would cause softmax to converge to one-hot distributions, reducing model expressiveness.

...

Index

agency, [9](#)

agent

 definition, [9](#)

attention mechanism, [9](#)

autonomy, [9](#)

ELIZA, [9](#)

GloVe, [9](#)

Large Language Models, [9](#)

license, [2](#)

self-supervised learning, [10](#)

tokenization, [10](#)

Transformer, [9](#)

Word2Vec, [9](#)