# Index

**A**

accuracy
confusion matrices, [Using a Task-Specific Model](#)
output verification, [Output Verification](#)
adaptive pretraining, [Using TSDAE for Domain Adaptation](#)
agents, [Agents: Creating a System of LLMs](#)-[ReAct in LangChain](#)
agentic RAG, [Agentic RAG](#)
ReAct in LangChain, [ReAct in LangChain](#)-[ReAct in LangChain](#)
step-by-step reasoning, [The Driving Power Behind Agents: Step-by-step Reasoning](#)-[The Driving Power Behind Agents: Step-by-step Reasoning](#)
AI (artificial intelligence)
accelerated development of, [An Introduction to Large Language Models](#)
defined, [What Is Language AI?](#)
ALBERT, [Model Selection](#)
align_labels function, [Preparing Data for Named-Entity Recognition](#)
all-MiniLM-L6-v2 model, [Supervised](#)
all-mpnet-base-v2 model, [Fine-Tuning for Few-Shot Classification](#)
Annoy, [Nearest neighbor search versus vector databases](#)
Anthropic Claude, [Proprietary, Private Models](#)
APIs (application programming interfaces), [Proprietary, Private Models](#)
Cohere, [API Keys](#), [Dense retrieval example](#)
external, [ChatGPT for Classification](#)
generating embeddings, [Supervised Classification](#)
OpenAI, [API Keys](#), [ChatGPT for Classification](#)
artificial intelligence (see AI)
ArXiv, [ArXiv's Articles: Computation and Language](#)

**B**

**M**

## Q

**U**

**V**