# Index

## About the Authors

**Jay Alammar** is Director and Engineering Fellow at Cohere (pioneering provider of large language models as an API). In this role, he advises and educates enterprises and the developer community on using language models for practical use cases. Through his popular AI/ML blog, Jay has helped millions of researchers and engineers visually understand machine learning tools and concepts from the basic (ending up in the documentation of packages like NumPy and pandas) to the cutting-edge (Transformers, BERT, GPT-3, Stable Diffusion). Jay is also a co-creator of popular machine learning and natural language processing courses on Deeplearning.ai and Udacity.

**Maarten Grootendorst** is a Senior Clinical Data Scientist at IKNL (Netherlands Comprehensive Cancer Organization). He holds master's degrees in organizational psychology, clinical psychology, and data science, which he leverages to communicate complex machine learning concepts to a wide audience. With his popular blogs, he has reached millions of readers by explaining the fundamentals of artificial intelligence—often from a psychological point of view. He is the author and maintainer of several open source packages that rely on the strength of large language models, such as BERTopic, PolyFuzz, and KeyBERT. His packages are downloaded millions of times and used by data professionals and organizations worldwide.

## Colophon

The animal on the cover of *Hands-On Large Language Models* is a red kangaroo (*Osphranter rufus*). They are the largest of all kangaroos, with a body length that can get up to a little over 5 feet and a tail as long as 3 feet. They are very fast and can hop to speeds over 35 miles per hour. They can jump 6 feet high and leap a distance of 25 feet in a single bound. The position of their eyes allows them see up to 300 degrees.

Red kangaroos are named after the color of their fur. While the name makes sense for the males—they have short, red-brown fur—females are typically more of a blue-grey color with a tinge of brown throughout. The red color in their fur comes from a red oil excreted from the glands in their skin. Because of their color, Australians refer to male red kangaroos as "big reds." However, because females are faster than males, they are often called "blue fliers."

Preferring open, dry areas with some trees for shade, red kangaroos can be found across Australia's mainland except in the upper north, lower southwest, and east coast regions of the country. Surrounding environmental conditions can affect reproduction. Because of this, females can pause or postpone pregnancy or birth until conditions are better. They often use this ability to delay birth of a new baby (joey) until the previous one has left their pouch.

The cover illustration is by Karen Montgomery, based on an antique line engraving from *Cassell's Popular Natural History*. The series design is by Edie Freedman, Ellie Volckhausen, and Karen Montgomery. The cover fonts are Gilroy Semibold and Guardian Sans. The text font is Adobe Minion Pro; the heading font is Adobe Myriad Condensed; and the code font is Dalton Maag's Ubuntu Mono.

# O'REILLY®

# Learn from experts.
# Become one yourself.

Books | Live online courses
Instant answers | Virtual events
Videos | Interactive learning

Get started at oreilly.com.