

## **Hard PDF Text for Extraction (Hyphenation + Line Breaks)**

### **English (with hyphenated line breaks)**

This document is designed to stress-test PDF text extraction. It contains hyphenation at line breaks, plus irregular spacing and newlines.

The course structure was generally coherent, but the assignments felt underspecified. Some students reported that feedback arrived too late, while others found it exceptionally helpful.

In week three we covered tokenization, normalization, and sentence segmentation. The lecture notes were well-written, yet a few sections were needlessly verbose. The project was motivating, but the evaluation rubric was unclear.

A more realistic example (multiple line breaks):

We learned about embeddings  
and vector retrieval,  
then built a small RAG pipeline to reduce hallucinations.

### **Danish (med trukne**

rede ord og linje-skift)

Dette dokument er lavet til at teste udtræk af tekst fra PDF. Det indeholder ord, der bliver delt ved linje-skift, og det kan ødelegge tokenisering, hvis man ikke ryder op først.

Kurset var overordnet set godt struktureret, men nogle krav var uklare. Flere studerende savnede hurtig feedback, især tidligt i forløbet, og det gjorde planlægning frustrende.

I uge tre arbejdede vi med forbehandling: tokenisering, normalisering og sætningssegmentering. Underviseren var dygtig til at forklare, men tempoet var til tider ujævnt.

Et "paper-agtigt" eksempel med ekstra linje-skift:

Vi talte om embeddings  
og vektor-søgning,  
og bagefter byggede vi en lille RAG-pipeline for at mindske hallucinationer.