

PDF of scanned
document



PNG images of
each page



First pass OCR



Prepare
postcorrection
dataset



Manual correction
and alignment