# Photo OCR

- Problem description and pipeline

    photo OCR : photo Optical Character Recognition

- Photo OCR Pipeline

    1. Text detection
    2. character segmentation
    3. character classification

    Each above steps may be a machine learning problem

    Pipeline : A system with many components / stages, several of which may use machine learning

- sliding Windows

    Text detection.

# Getting lots of data : Artificial data synthesis

- Creating new data from scratch.

    Take real data and add background / noise / distortion

    Distortion should represent the type of noise in the test set.

- Discussion on getting more data.

    1. Make sure you have low bias classifier before expanding the effort.
    2. How much work would it be to get 10X much data as we

currently have?

- Artificial data synthesis?
- Collect/label it yourself?
- Crowd source (e.g. Amazon mechanical Turk

Ceiling Analysis: What part of pipeline to work on next?

Image $\longrightarrow$ Text detection $\longrightarrow$ character segmentation $\longrightarrow$ character recognition

Ceiling analysis: estimating the error due to each component