# GETTING GOOD DATA* FASTER

*) human labeled data sets
+ machine learning solutions

## MACHINE LEARNING DEVELOPMENT CYCLE ↻

training → evaluation → "patching" production results

Sub cycle:

## HLD PRODUCTION PROCES
define → pilot → produce

## LABELING WORKFORCES
"How core is this to my business?"

## PRODUCTION ARCHITECTURE
aka dataset pipeline