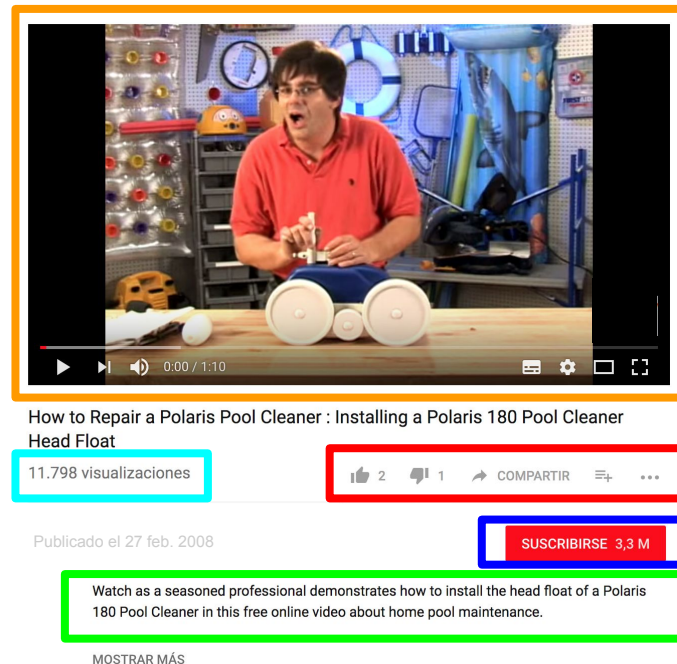# Dataset

- 2000h of **how-to** videos (Yu et al., 2014)
  - 300h for MT, 480h for ASR (as of today)
  - Shared splits, held-out data
- Ground truth captions
- Metadata
  - Number of likes / dislikes
  - Visualizations
  - Uploader, Date
  - Tags
- Video descriptions ("summaries")
  - 80K descriptions for 2000h
- Very different topics
  - Cooking, fixing things, playing instruments, etc.
- 300,000 segments translated into Portuguese



How to Repair a Polaris Pool Cleaner : Installing a Polaris 180 Pool Cleaner Head Float

11.798 visualizaciones

👍 2   👎 1   ➦ COMPARTIR   ⊟₊   •••

Publicado el 27 feb. 2008

SUSCRIBIRSE  3,3 M

Watch as a seasoned professional demonstrates how to install the head float of a Polaris 180 Pool Cleaner in this free online video about home pool maintenance.

MOSTRAR MÁS