# Evaluation Campaigns

1. Speech
2. Others
3. General Points

Dan Ellis

Dept. Electrical Engineering, Columbia University
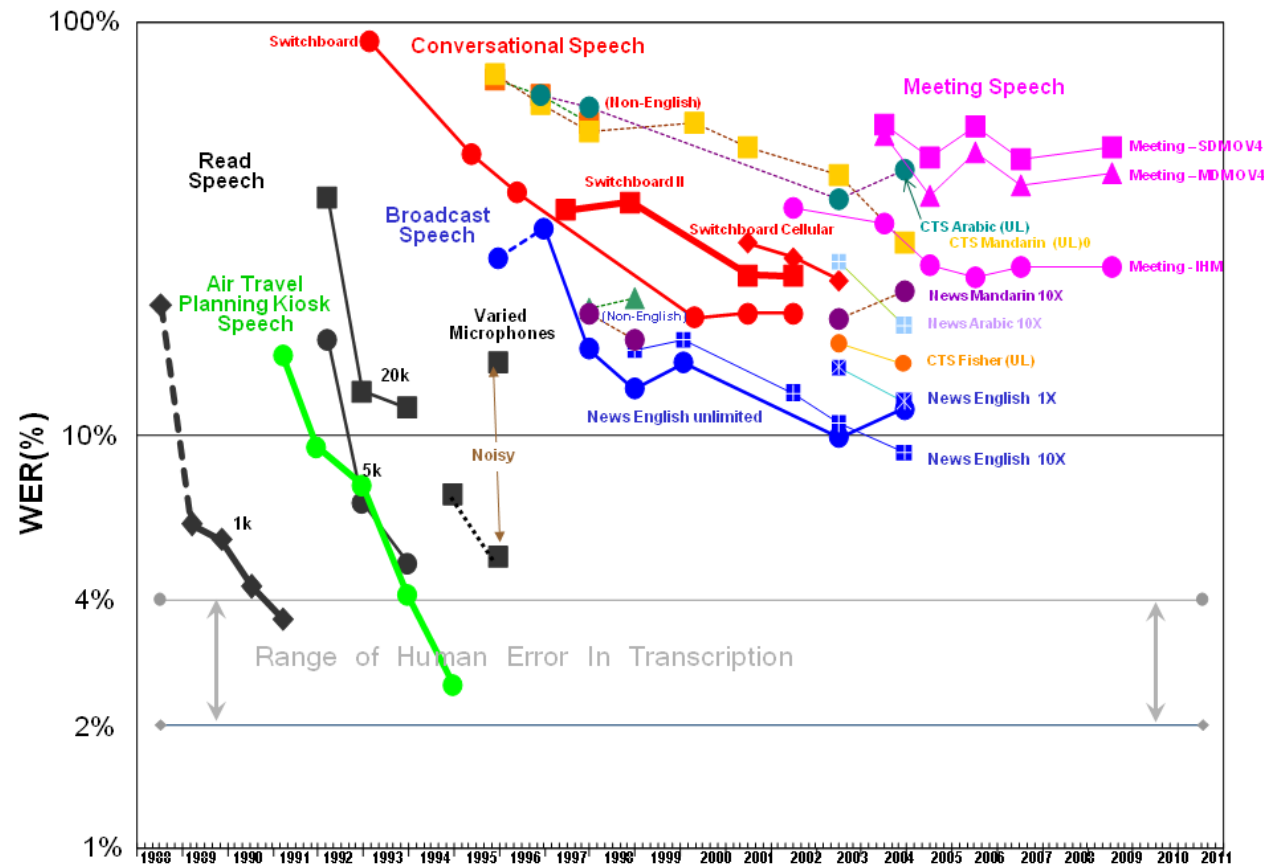
dpwe@ee.columbia.edu    http://www.ee.columbia.edu/~dpwe/e4896/

# Evaluations

- Systematically evaluating research output with common data & metrics

- DARPA/NIST Speech Recognition as the original & canonical example:



NIST STT Benchmark Test History – May. '09

http://www.itl.nist.gov/iad/mig/publications/ASRhistory/

# The Origin of Evaluations
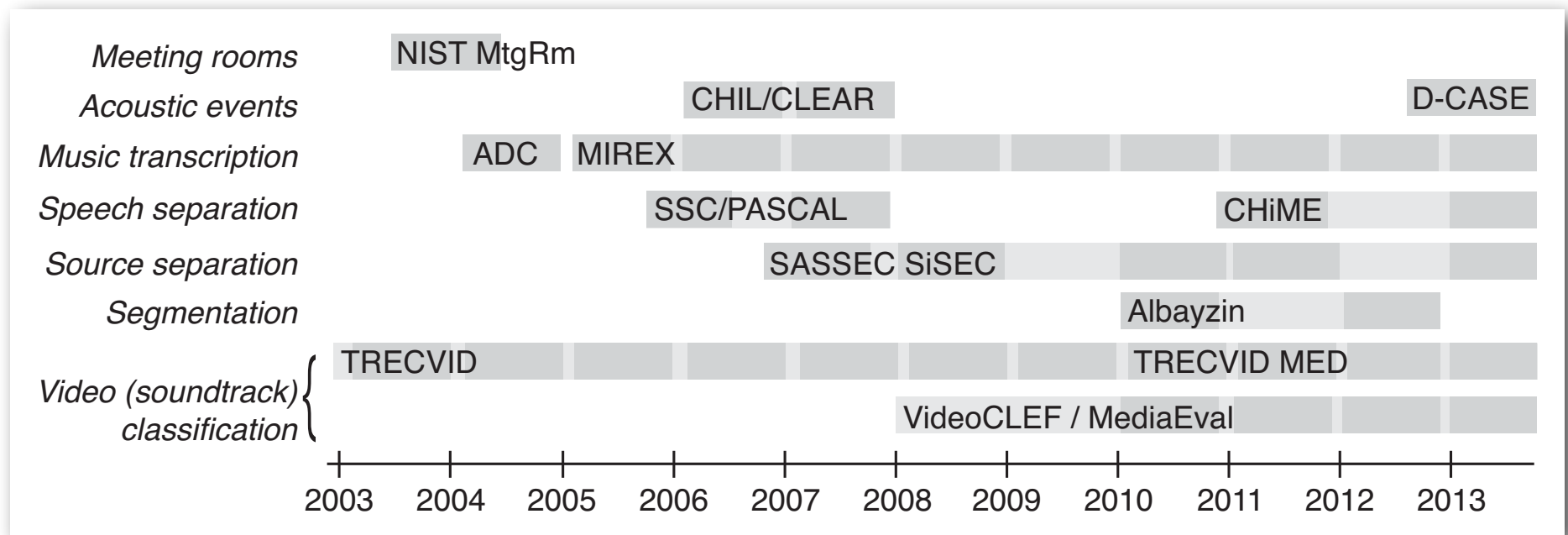


- Mark Liberman: "Avoiding glamour and deceit"

  ○ placate funders!

- American Association for the Advancement of Science Meeting, 2011-02-19,
  The Digitization of Science: Reproducibility and Interdisciplinary Knowledge Transfer
  http://languagelog.ldc.upenn.edu/nll/?p=2976
  http://www.stanford.edu/~vcs/AAAS2011/AAAS2011Liberman.pdf

# Other Evaluations

- **Benefits to speech led to many copies**
  - now standard for DARPA and IARPA programs
  - emulated in many other fields
  - typically volunteer-funded

- **Example: Sound Scene Analysis evaluations**



| | 2003 | 2004 | 2005 | 2006 | 2007 | 2008 | 2009 | 2010 | 2011 | 2012 | 2013 |
|---|---|---|---|---|---|---|---|---|---|---|---|
| *Meeting rooms* | NIST MtgRm | | | | | | | | | | |
| *Acoustic events* | | | | CHIL/CLEAR | | | | | | | D-CASE |
| *Music transcription* | ADC | MIREX | | | | | | | | | |
| *Speech separation* | | | SSC/PASCAL | | | | | CHiME | | | |
| *Source separation* | | | | SASSEC SiSEC | | | | | | | |
| *Segmentation* | | | | | | | Albayzin | | | | |
| *Video (soundtrack) classification* | TRECVID | | | | | | TRECVID MED | | | | |
| | | | | | | | VideoCLEF / MediaEval | | | | |

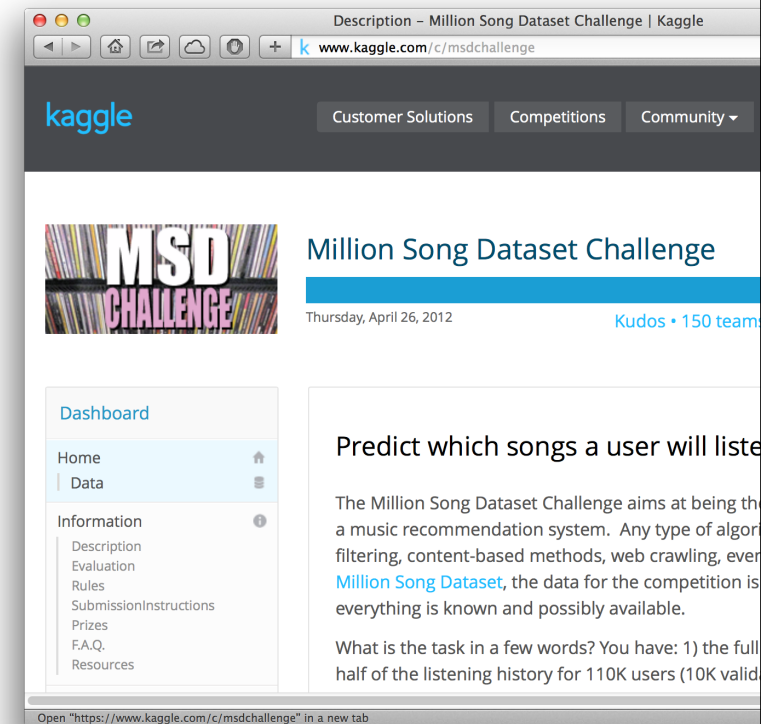  - Metrics:  SNR,  Frame Acc,  Event Error Rate,  mAP

# MIREX

- **Effort to compare Music Information Retrieval algorithms**
  - ○ organized by Stephen Downie, UIUC
  - ○ funded by Mellon Foundation

- **First round in 2005**
  - ○ 5-15 tasks per year, 3-20 participants per trial
    = 2037 algorithms run over 37 datasets in 10 years

- **Impact**
  - ○ organized, solidified research areas - chords, covers
    - focus of community discussion of agenda
  - ○ no open release of data - you have to participate

# Million Song Dataset Challenge

- **Listening history data for 1M+ listeners**
  - but not time-stamped
  - task is to rank tracks based on partial history

- **kaggle.com:** "predictive analytics leader"
  - actually, a platform for big-data challenges

- **Competition**
  - ran for 4 months in 2012; 150 teams participated
  - avg. prec. improved from 0.024 to 0.179
  - .. but no audio features used!

# Aspects of Evaluations

- Relevance of task & metrics
  - at least you'll solve one task

- Scale matters
  - for statistical significance & non-over-fitting

- Encouraging participation
  - plusses and minuses of participating

- Models for distributing the effort
  - it's a lot of work to run these systems; who pays? +secrets

- Ensuring the sharing of information
  - opportunity to share code?

- Releasing test materials
  - .. for extensive post-mortems … but next time?

# Impact of Evaluations

- Good:
  - direct comparison of techniques
    - invest with confidence!
  - focus community research effort

- Bad:
  - non-evaluated topics are starved of attention
  - leads to conservative monoculture
    - puts off good newcomers?
  - too much focus on one number…

# Summary

- **Glamour and Deceit**
  - common data & tasks provide clarity

- **Knowledge and Progress**
  - identify the things that work
    (and how they combine)

- **Data and Code**
  - systems that conform to a common standard
    are (more) ready for sharing