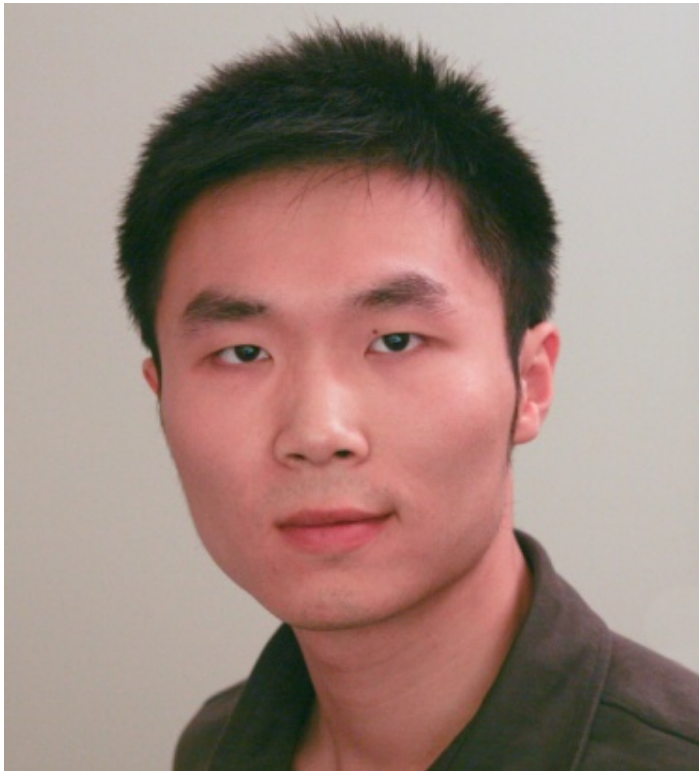# Large-Scale Visual Recognition Powered by Big Data and Big Crowd

Fei-Fei Li

Stanford University

Dr. Jia Deng
Stanford U. -> U. Michigan

Prof. Kai Li
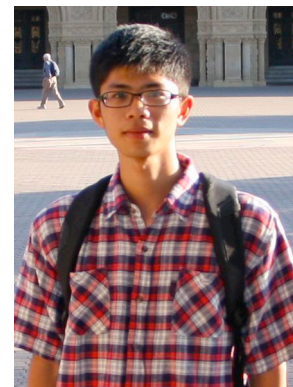Princeton U.

Prof. Alex Berg
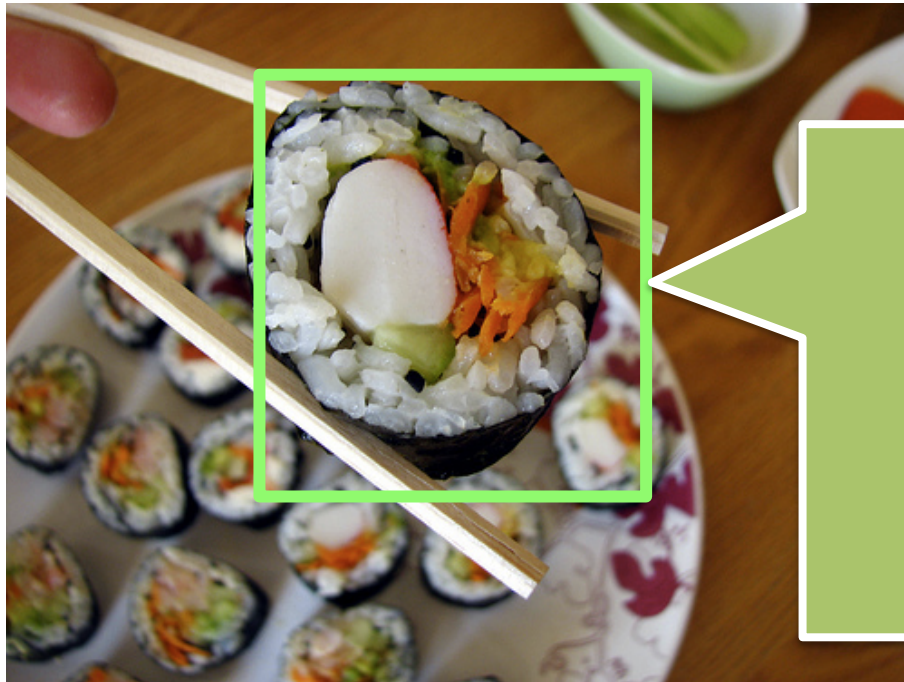Stony Brook U.

Sanjeev Satheesh
Stanford U.

Jonathan Krause
Stanford U.

Zhiheng Huang
Stanford U.

Olga Russakovsky
Stanford U.

**California Roll**

*Ingredients*: Rice, Seaweed, Crab, Cucumber, Avocado

*Calories*: 40

*Fat*: 7g

*Carb*: 40g

*Protein*: 5g

Gluten Free

**Amanita phalloides**
http://en.wikipedia.org/
wiki/Amanita_phalloides

TOXIC. DO NOT EAT

IKEA POANG Chair
ON SALE
$29.00 at ikea.com

**Mornonga
(Japanese flying squirrel)**
Inhabits sub-alpine forests in Japan.
Nocturnal. Eats seeds, fruit, tree leaves
(Wikipedia)

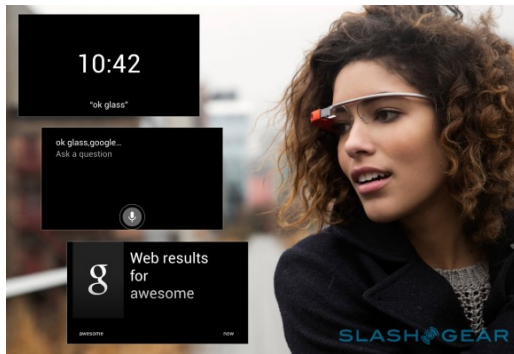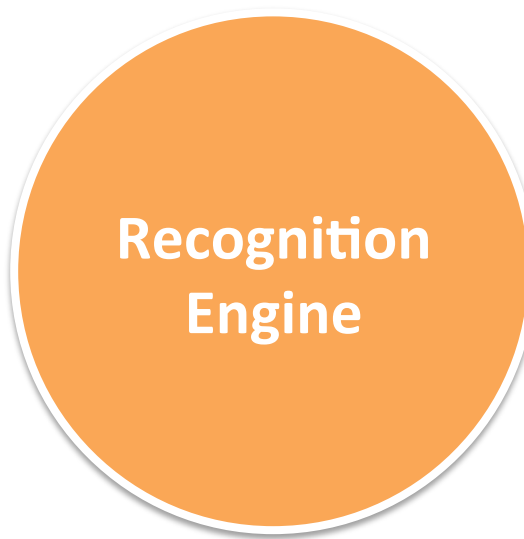**I wish my computer to recognize EVERYTHING**

Surveillance


Robotics


Assistive tools


Wearable devices

**Recognition Engine**
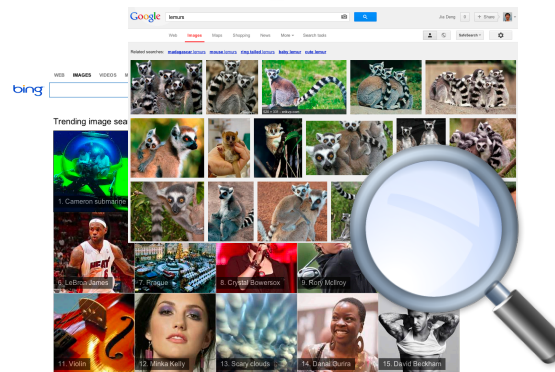

Driverless cars


Smart photo album


Image search


Mining social media

# What can computers already recognize?

The Nikon S60. Detects up to 12 faces.

# Google Goggles

Use pictures to search the web.

New!

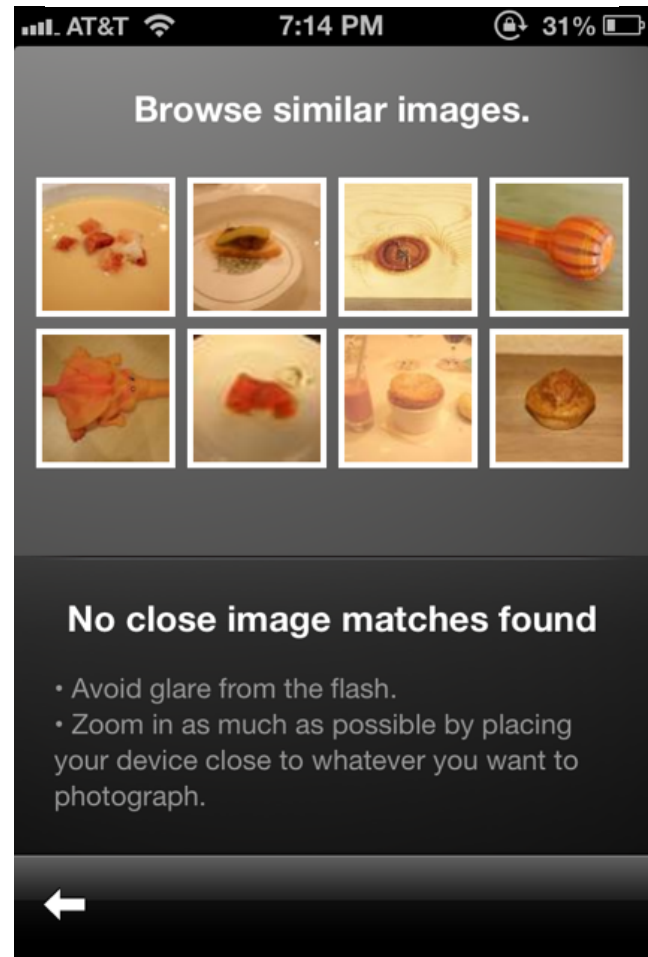| Text | Landmarks | Books | Contact Info | Artwork | Wine | Logos |

German (auto) » English

Lammkoteletts vom Biobauern mit Schalotten, Tomatencoulis und Basilikum-Gnocchi

Lamb chops from the farmers with the shallots, tomato sauce and basil gnocchi

.ull. AT&T 📶          7:14 PM          @ 31% 🔋

**Browse similar images.**

**No close image matches found**

• Avoid glare from the flash.
• Zoom in as much as possible by placing your device close to whatever you want to photograph.

# PASCAL VOC [Everingham et al. 2006-2012]



| | |
|---|---|
| **Airplane** | **Dining table** |
| **Bird** | **Dog** |
| **Boat** | **Horse** |
| **Bike** | **Motorbike** |
| **Bottle** | **Person** |
| **Bus** | **Potted plant** |
| **Car** | **Sheep** |
| **Cat** | **Sofa** |
| **Chair** | **Train** |
| **Cow** | **TV monitor** |

# No Coffee Mugs!

# What about Gas Pumps!

# Many Things Do Not Work Yet…

# How many things are there?

WordNet

**10K+**
[Biederman '87]

**60K+**
product
categories

**80K+**
English nouns
[Miller '95; Fellbaum '98]

**3.5M+**
unique tags
[Sigurbjörnsson & Zwol '08]

**4.1M+**
articles

From 20 classes to Millions?

# nature

**THE BITER BIT**
Viral infections for viruses

**TROPICAL CYCLONES**
The strong get stronger

**BLACK HOLE PHYSICS**
A new window on the
Galactic Centre

BIG DATA

**NATUREJOBS**
Minnesota musings

# SCIENCE IN THE PETABYTE ERA

# Big Data from the Internet

# Global Consumer Internet Traffic Per Month



| Year | Traffic |
|------|---------|
| 2011 | 21 EB |
| 2012 | 30 EB |
| 2013 | 38 EB |
| 2014 | 48 EB |
| 2015 | 63 EB |
| 2016 | 83 EB |

Source: Cisco

**YouTube** 72 hours of videos / min

**facebook** 300 million images / day

83 EB

Visual

86%

2016

Source: Cisco

# Big Data from the Internet

# → The Internet can teach EVERYTHING

Google pitbullfrog

**Evolution Gone Wild**
*Future plants and animals*

http://www.worth1000.com/contests/12705/contest

The Internet: Machines + Crowd

Big Data

**Teach machines to recognize EVERYTHING**

**PASCAL VOC**

20
[Everingham et al.'06-'12]

10K+
[Biederman '87]

# Goal: Build a recognition engine on 10K classes

## And Never Make A Mistake!

**EVA**
*Engine for Visual Annotation*

**The EVA system**, powered by **ImageNet**, can annotate images with guaranteed accuracies. It currently recognizes over **10,000** visual categories. See the project page to find out more.

**Paste a URL** | Upload an image

ANNOTATE

# Agenda

**How to build a large-scale recognition engine using big data**

**STEP 1:** ?

**STEP 2:** ?

**STEP 3:** ?

# Agenda

**How to build a large-scale recognition engine using big data**

**STEP 1:** | **Build a Large Knowledge Base**

**STEP 2:** | **?**

**STEP 3:** | **?**

**Get a list of everything**

WordNet

80K nouns

- Expert constructed
- Rich structure
  - Taxonomy, Partonomy
- Widely used

**Crawl the web**

Google flickr

[Torralba, Fergus, Freeman '08]
[Yao, Yang, Zhu '07]
[Everingham et al '06]
[Russell et al '05]
[Griffin & Perona '03]
[Fei-Fei, Fergus, Perona '03]

# What you get may not be what you want.



[Torralba, Fergus, Freeman '08]

bluebird

bluebird

**Get a list of everything**

WordNet

80K nouns

- Expert constructed
- Rich structure
  - Taxonomy, Partonomy
- Senses disambiguated
- Widely used

**Crawl the web**

Google flickr

[Torralba, Fergus, Freeman '08]
[Yao, Yang, Zhu '07]
[Everingham et al '06]
[Russell et al '05]
[Griffin & Perona '03]
[Fei-Fei, Fergus, Perona '03]

**Clean up**

# Graduate Students

Good at complex tasks

Good quality

Very few of them

High cost

**Estimate: 20 Years, $2M+**

# Graduate Students

# The Crowd

Good at complex tasks

Good quality

Very few of them

High cost

# Graduate Students

Good at complex tasks

Good quality

Very few of them

High cost

# The Crowd

Good at simple tasks

Mixed quality

Many of them

Low cost

# Cleaning up by AMT workers

# Cleaning up by AMT workers

# Dealing with Mixed Quality



Quality Control
(Probabilistic
Models)

More?

N

Y

[Sheng et al. '08]
[Sorokin & Forsyth '08]
**[Deng et al. '09]**

# Bluebird

Blue North American songbird

ⓘ Numbers in brackets: (the number of synsets in the subtree ).

- ImageNet 2011 Fall Release (21841
  - animal, animate being, beast, bru
    - mate (0)
    - chordate (2953)
      - tunicate, urochordate, uroc
      - cephalochordate (1)
      - vertebrate, craniate (2943)
        - mammal, mammalian
        - aquatic vertebrate (578)
        - tetrapod (1)
        - amniote (0)
        - fetus, foetus (2)
        - Amniota (0)
        - amphibian (93)
        - reptile, reptilian (267)
        - bird (855)
          - dickeybird, dickey-bi
          - nonpasserine bird (
          - bird of prey, raptor, r
          - gallinaceous bird, g
          - parrot (19)
          - cuculiform bird (8)
          - coraciiform bird (14)
          - apodiform bird (8)
          - caprimulgiform bird
          - piciform bird (20)
          - trogon (2)
          - aquatic bird (278)
          - passerine, passerif
            - wren, jenny wren

| Treemap Visualization | Images of the Synset | Downloads |



*Images of children synsets are not included. All images shown are thumbnails. Images may be subject to copyright.

Prev | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | ... | 35 | 36 | Next

# IM**A**GENET [Deng et al. 2009]

**www.image-net.org**

**22K** categories and **14M** images

- Animals
  - Bird
  - Fish
  - Mammal
  - Invertebrate
- Plants
  - Tree
  - Flower
  - Food
  - Materials
- Structures
- Artifact
  - Tools
  - Appliances
  - Structures
- Person
- Scenes
  - Indoor
  - Geological Formations
- Sport Activities

# Number of Labeled Images

**ImageNet, 14M**

**[Deng et al. '09]**

SUN, **131K**

[Xiao et al. '10]

LabelMe, **37K**

[Russell et al. '07]

PASCAL VOC, **30K**

[Everingham et al. '06-'12]

Caltech101, **9K**

[Fei-Fei, Fergus, Perona, '03]

IMAGENET hired **50K+** AMT workers

who looked at **160M+** images

and made **550M+** binary decisions

U.S. economy outlook (Gallup)
Number of images in ImageNet

14M
12M
11M
10M
3M
0M

Jan-08  May-08  Sep-08  Jan-09  May-09  Sep-09  Jan-10  May-10  Sep-10  Jan-11  May-11

# Research Impact of ImageNet

**3000+** registered users, visits from 175 countries



1       75,394

ECCV 2012
**Best paper Award**

Kuettel, Guillaumin, Ferrari. **Segmentation Propagation in ImageNet.** ECCV 2012

Le et al. **Building high-level features using large scale unsupervised learning**. ICML 2012.

Krizhevsky, Sutskever, Hinton. **ImageNet classification with deep convolutional neural networks.** NIPS 2012

The New York Times

# Science

Search All NYTimes.com

Orange Savings Account℠

WORLD | U.S. | N.Y. / REGION | BUSINESS | TECHNOLOGY | SCIENCE | HEALTH | SPORTS | OPINION | ARTS | STYLE | TRAVEL | JOBS | REAL ESTATE | AUTOS

## Seeking a Better Way to Find Web Images

By JOHN MARKOFF
Published: November 19, 2012

STANFORD, Calif. — You may think you can find almost anything on the Internet.

**Connect With Us on Social Media**
@nytimesscience on Twitter.

· Science Reporters and Editors on Twitter

Like the science desk on Facebook.

But even as images and video rapidly come to dominate the Web, search engines can ordinarily find a given image only if the text entered by a searcher matches the text with which it was labeled. And the labels can be unreliable, unhelpful ("fuzzy" instead of "rabbit") or simply nonexistent.

To eliminate those limits, scientists will need to create a new generation of visual search technologies — or else, as the Stanford computer scientist Fei-Fei Li recently put it, the Web will be in danger of "going dark."

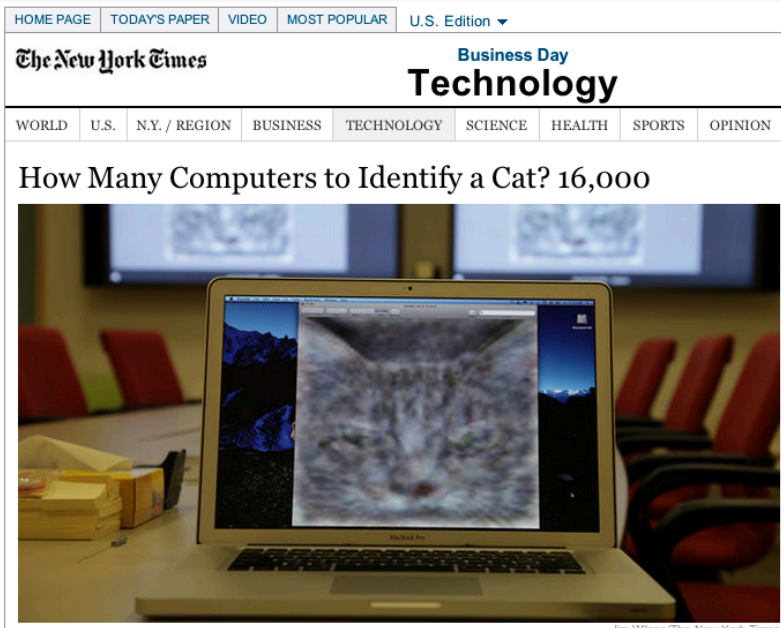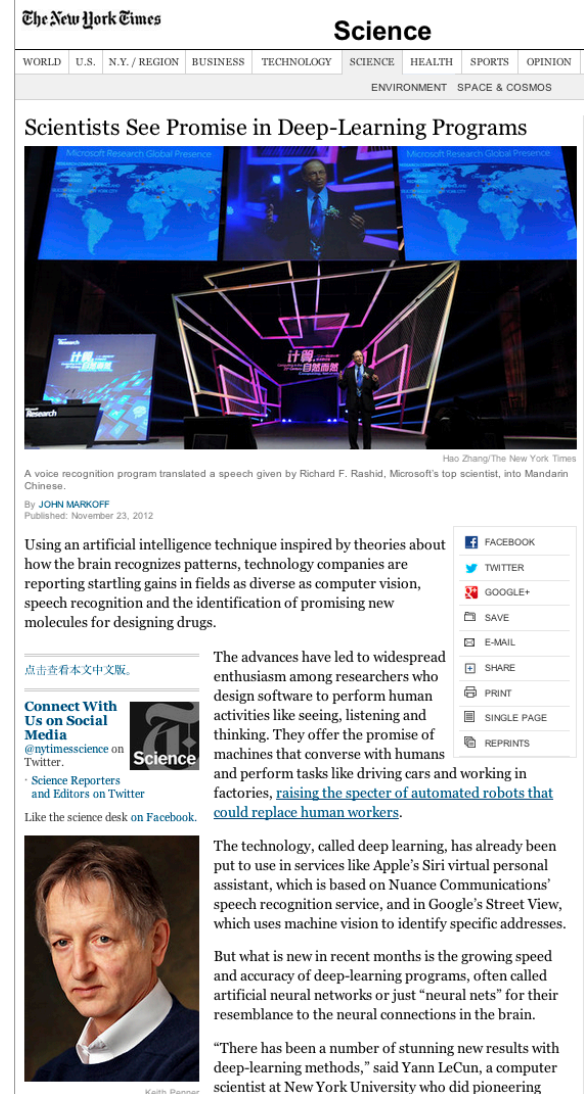Now, along with computer scientists from Princeton, Dr. Li, 36, has built the world's largest visual database in an effort to mimic the human vision system. With more than 14 million labeled objects, from obsidian to orangutans to ocelots, the database has because a vital resource for computer vision researchers.

The labels were created by humans. But now machines can learn from the vast database to recognize similar, unlabeled objects, making possible a striking increase in recognition accuracy.

This summer, for example, two Google computer scientists, Andrew Y. Ng and Jeff Dean, tested the new system, known as ImageNet, on a huge collection of labeled photos.

FACEBOOK

TWITTER

GOOGLE+

SAVE

E-MAIL

SHARE

PRINT

REPRINTS

THE SESSIONS
NOW PLAYING

# Agenda

**How to build a large-scale recognition engine using big data**

STEP 1: | **Build a Large Knowledge Base (ImageNet)**

STEP 2: | **?**

STEP 3: | **?**

# Learn to Classify 10K Classes

- 9 Million images

- 4 methods
  - SPM+SVM [Lazebnik et al. '06]
  - BOW+SVM [Csurka et al. '04]
  - BOW+NN
  - GIST+NN [Oliva et al. '01]

- 6.4% for 10K categories



Deng, Berg, Li, & Fei-Fei, *ECCV2010*

# Learn to Classify 10K Classes



Deng, Berg, Li, & Fei-Fei, *ECCV2010*

# Fine-grained categories are a lot harder



Average Semantic Distance

Finer

Coarser

Deng, Berg, Li, & Fei-Fei, *ECCV2010*

# Agenda

**How to build a large-scale recognition engine using big data**

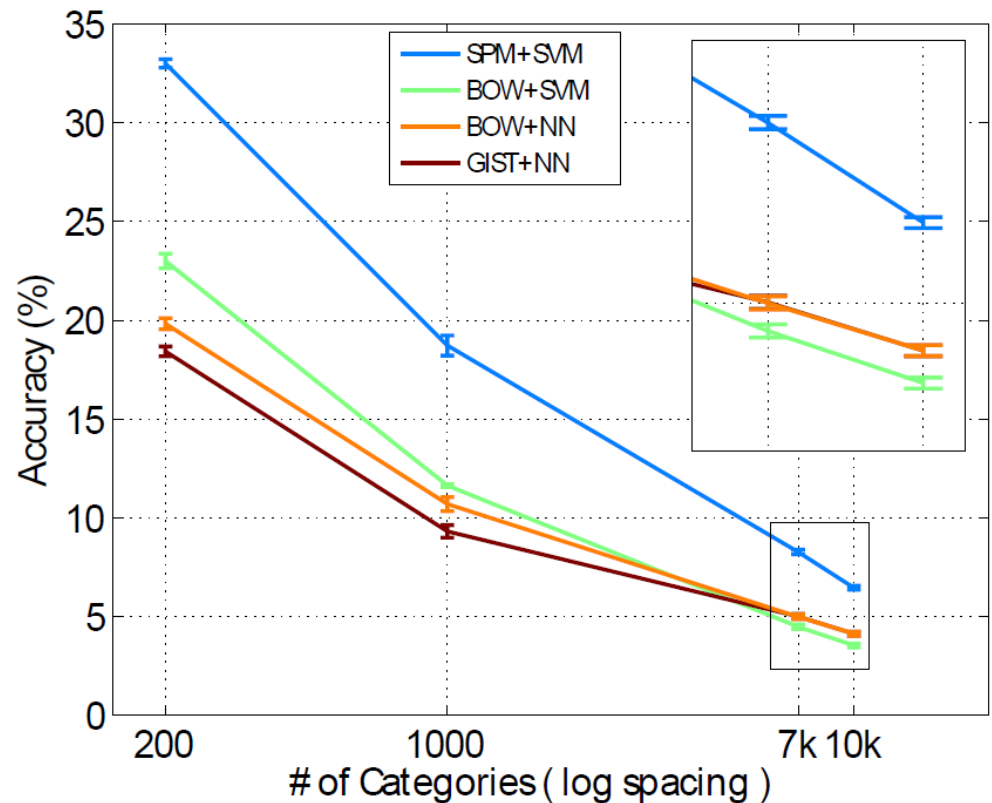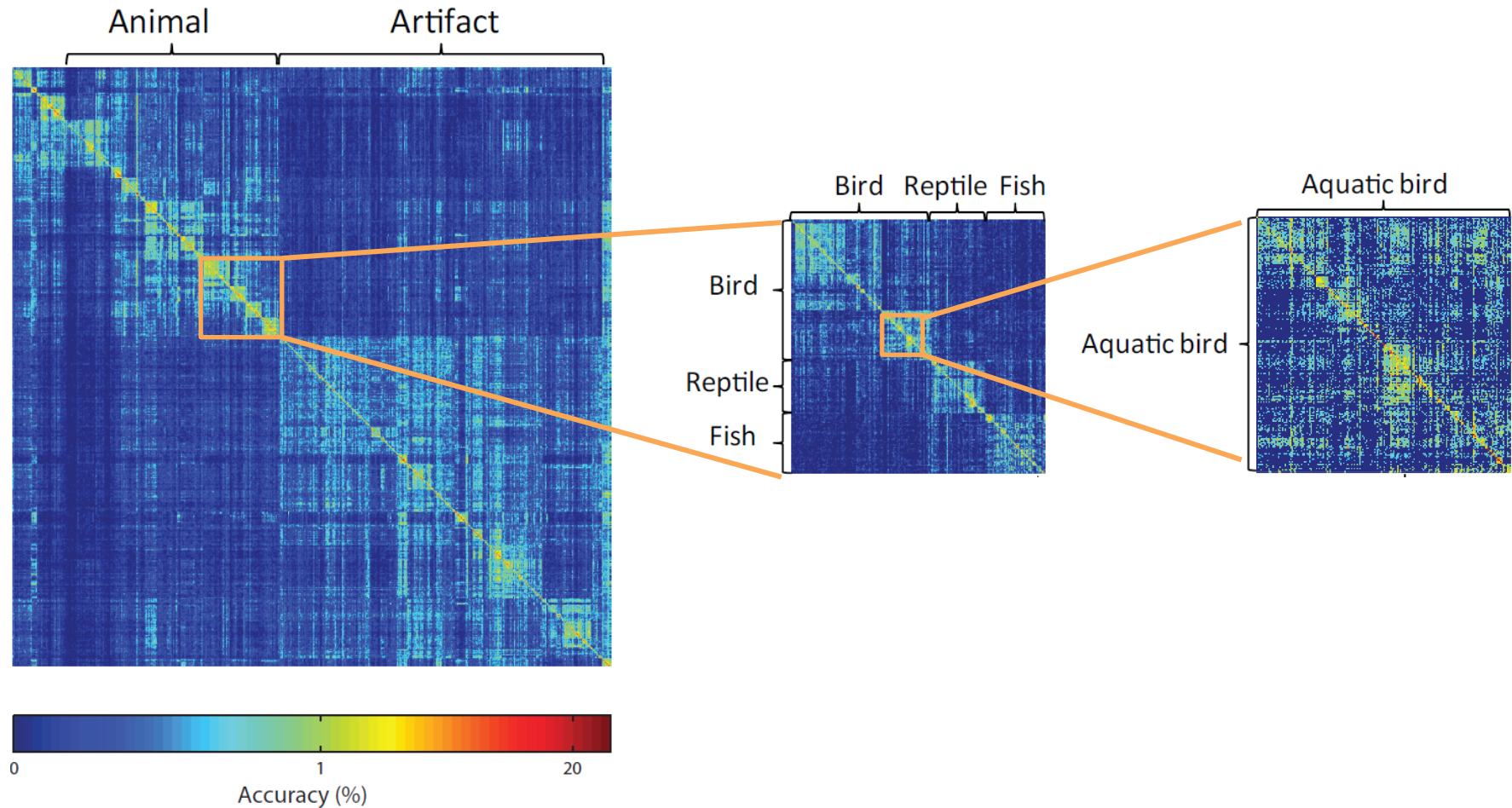**STEP 1:** Build a Large Knowledge Base (ImageNet)

**STEP 2:** Fine-Grained Recognition

**STEP 3:** ?

# Why is Fine-Grained Recognition Difficult?



**What breed is this dog?**

# Why is Fine-Grained Recognition Difficult?



Cardigan Welsh Corgi

Pembroke Welsh Corgi

What breed is this dog?

**Key: Find the right features.**
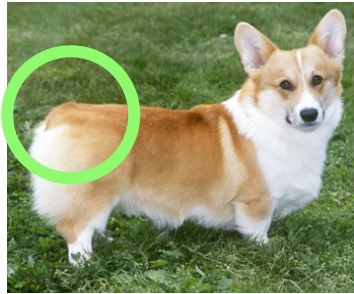
# Why is Fine-Grained Recognition Difficult?



rdigan Welsh Corgi

mbroke Welsh Corg

**Learning**

**Existing Work**

[Branson et al. '10]

[Bo et al. '10]

[Farrell et al. '11]

[Yao et al. '11]

[Yao et al. '12]

# Why is Fine-Grained Recognition Difficult?



Cardigan Welsh Corgi

Pembroke Welsh Corgi

**Learning**

**Existing Work**

[Branson et al. '10]

[Bo et al. '10]

[Farrell et al. '11]

[Yao et al. '11]

[Yao et al. '12]

# Why is Fine-Grained Recognition Difficult?



Cardigan Welsh Corgi

Pembroke Welsh Corgi

Learning

**How to help computers select features?**

# Machine-Crowd Collaboration

**Machine**

**KNOWLEDGE**

**Crowd**

# Machine-Crowd Collaboration

# Machine-Crowd Collaboration

[Prairie Warbler (wikipedia)](wikipedia)

[Yellow Warbler (wikipedia)](wikipedia)

# Machine-Crowd Collaboration

# The BubbleBank Representation



Deng, Krause, & Fei-Fei, *CVPR2013*

# mAP on CUB-14 [Welinder et al. 10]



Bar chart values:
- MKL: 37.02
- Birdlet: 40.05
- CFAF: 44.73
- BubbleBank: 58.47

MKL [Branson et al. '10]
Birdlet [Farrell et al. '11]
CFAF [ Yao et al.'12]

# Accuracy on CUB-200 [Welinder et al. 10]



Bar chart values:
- LLC: 18
- MKL: 19
- RF: 19.2
- MultiCue: 22.4
- KDES: 26.2
- Tricos: 26.7
- BubbleBank: 32.8

MKL [Branson et al. '10]
LLC [Wang et al. '09]
RF [Yao et al. '11]
MultiCue [Khan et al.'11]
KDES [Bo et al. '10]
Tricos [Chai '12]

Deng, Krause, & Fei-Fei, *CVPR2013*

# Top Activated Bubbles (successful predictions)



Deng, Krause, & Fei-Fei, *CVPR2013*

# Agenda

**How to build a large-scale recognition engine using big data**

STEP 1: | **Build a Large Knowledge Base (ImageNet)**

STEP 2: | **Fine-Grained Recognition (Bubbles)**

STEP 3: | **?**

# Agenda

**How to build a large-scale recognition engine using big data**

STEP 1: **Build a Large Knowledge Base (ImageNet)**

STEP 2: **Fine-Grained Recognition (Bubbles)**

STEP 3: **Putting a label on "everything"**

# The Current State of the Art

| 10K classes | 32.6% | Krizhevsky et al. NIPS 2012 |
|---|---|---|
| 20K classes | 15% | Le et al. NIPS 2012 |

## Not quite practical yet…

## But we are measuring the very fine-grained level

# Hedging: Be as informative as possible with few mistakes



Deng, Krause, Berg, Fei-Fei, CVPR2012

# Formal Problem Statement



Deng, Krause, Berg, Fei-Fei, CVPR2012

# Formal Problem Statement



Deng, Krause, Berg, Fei-Fei, CVPR2012

# Formal Problem Statement



Deng, Krause, Berg, Fei-Fei, CVPR2012

# Formal Problem Statement

## Assumptions

- Same distribution for training and test.
- A base classifier **g** that gives posterior probability on the hierarchy.

## Goal

- Find a **decision rule f**
  - Expected accuracy **A(f)** is at least **1-ε**
  - Maximize expected reward **R(f)**

$$\underset{f}{Maximize} \quad R(f)$$

$$Subject \ to \quad A(f) \geq 1 - \varepsilon$$



Test image

**g**

posterior for all nodes

**f**

Deng, Krause, Berg, Fei-Fei, CVPR2012

# Pick a global confidence threshold T=0.9 [Vailaya et al. '99]



100 images

1.0 $0

0.90 $1

0.6 $2

Another 100 images

$0 1.0

$1 0.90

0.80

$10

Reward = ($1 * 0.90 + $1 * 0.90) / 2 = $0.90
Accuracy = (0.90 + 0.90 ) / 2 = 0.90

Deng, Krause, Berg, Fei-Fei, CVPR2012

# Pick a global confidence threshold T=0.9 [Vailaya et al. '99]



100 images

1.0 — $0
0.90 — $1
0.6 — $2

Another 100 images

$0 — 1.0
$1 — 0.90
0.80 — $10

Reward = ($1 * 0.90 + $1 * 0.90) / 2 = $0.90
Accuracy = (0.90 + 0.90) / 2 = 0.90

100 images

T=0.95

1.0 — $0
0.90 — $1
0.6 — $2

Another 100 images

$0 — 1.0
$1 — 0.90
T=0.80
0.80 — $10

Reward = ($0 * 1.0 + $10 * 0.80) / 2 = $4
Accuracy = (1.0 + 0.80) / 2 = 0.90

Deng, Krause, Berg, Fei-Fei, CVPR2012

**We can optimize individual thresholds...**

**But actually we don't need to.**
**There is a simpler and provably optimal solution**

Deng, Krause, Berg, Fei-Fei, CVPR2012

# A global, fixed scalar parameter *λ≥0*



Test image

posterior for all nodes

Increase each reward by *λ*

*λ*

Expected rewards

Predict the best

Deng, Krause, Berg, Fei-Fei, CVPR2012

# The DARTS algorithm



Theorem: Under very mild conditions, this is optimal.

Deng, Krause, Berg, Fei-Fei, CVPR2012

ImageNet10K

Ours
LEAF-GT
MAX-REW
MAX-EXP

Deng, Krause, Berg, Fei-Fei, CVPR2012

**The EVA system**, powered by **ImageNet**, can annotate images with guaranteed accuracies. It currently recognizes over **10,000** visual categories. See the project page to find out more.

**Paste a URL** | Upload an image

ANNOTATE

## Google Goggles
### Use pictures to search the web.

**Browse similar images.**

**No close image matches found**

- Avoid glare from the flash.
- Zoom in as much as possible by placing your device close to whatever you want to photograph.

## EVA
### Engine for Visual Annotation

0.95 coffee mug
0.97 mug
0.99 drinking vessel

Google images

EVA
Engine for Visual Annotation

Image size:
401 × 604

No other sizes of this image found.

Visually similar images - Report images

0.87 face , gas pump, person
0.90 face , gas pump

0.75 artifact, crater, matter, vertebrate
0.77 crater, matter, vertebrate
0.78 chordate, crater, matter
0.86 animal, matter
0.87 animal

# Summary

**How to build a large-scale recognition engine using big data**

STEP 1: **Build a Large Knowledge Base (ImageNet)**

STEP 2: **Fine-Grained Recognition (Bubbles)**

STEP 3: **Putting a label on "everything" (Hedging)**

# ImageNet Challenge 2013

- 1.2 Million images and 1000 classes
- New PASCAL-style Detection Task
  - Full annotation of 200 classes in test images.
- http://www.image-net.org/challenges/LSVRC/2013/

# Fine-Grained Challenge 2013

- Competition on Fine-Grained Recognition
  - Airplanes, Birds, Cars, Dogs, Shoes.
- https://sites.google.com/site/fgcomp2013/
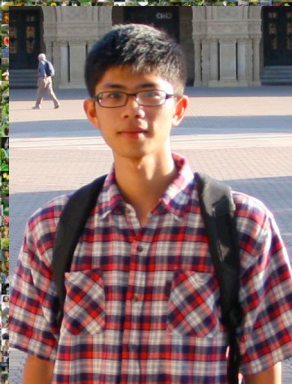
# Thank you!



Prof. Kai Li
Princeton U.

Prof. Alex Berg
Stony Brook U.

Sanjeev Satheesh
Stanford U.

Jonathan Krause
Stanford U.

Zhiheng Huang
Stanford U.

Olga Russakovsky
Stanford U.

Dr. Jia Deng
Stanford U.