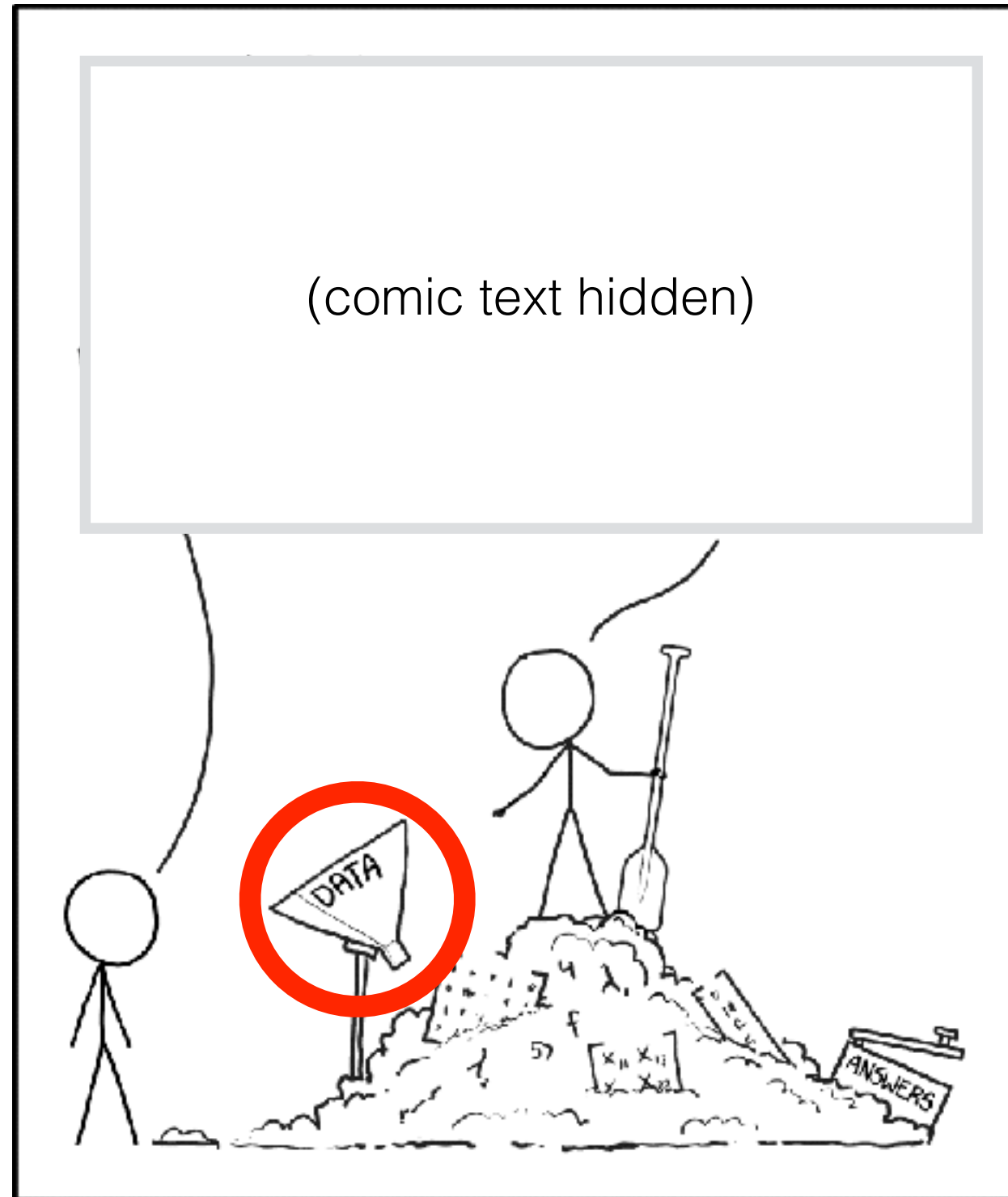


NEXT: Crowdsourcing, machine learning and cartoons

Scott Sievert
@stsievert  

Link to slides and proceedings:
tinyurl.com/scipy-next

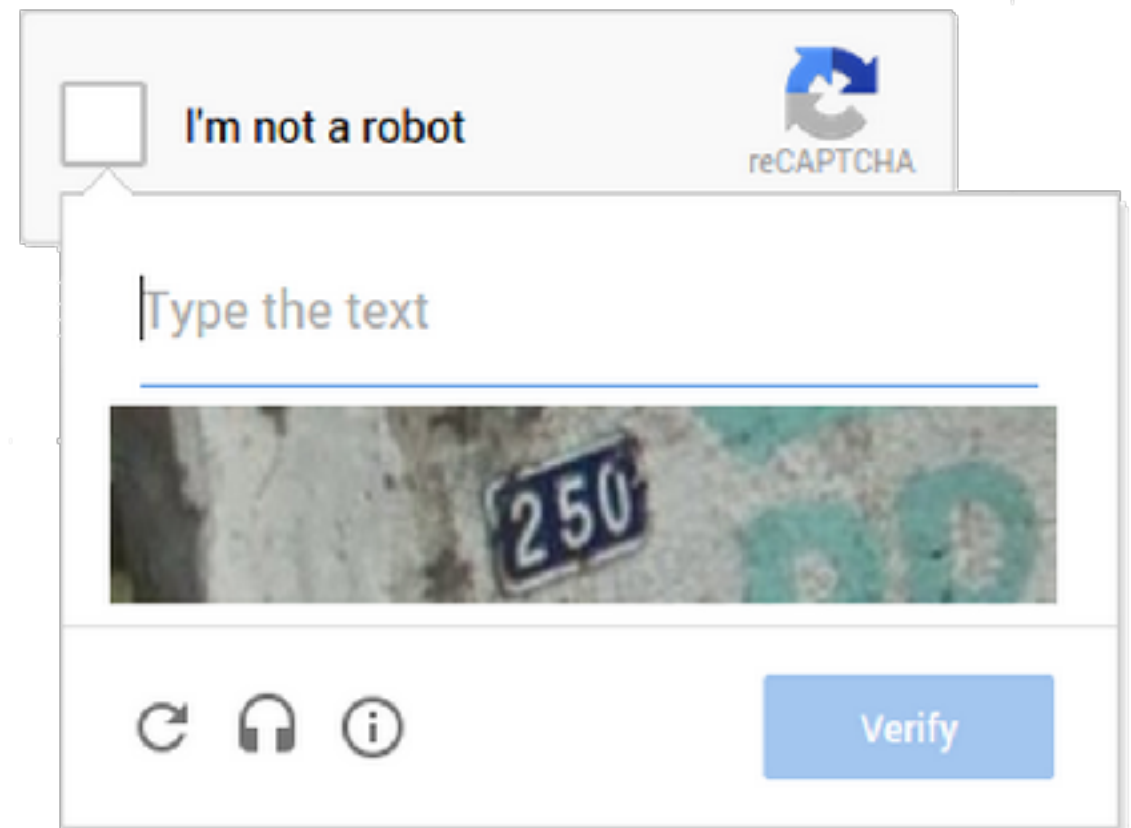
Problem



Data collection can be costly

Example

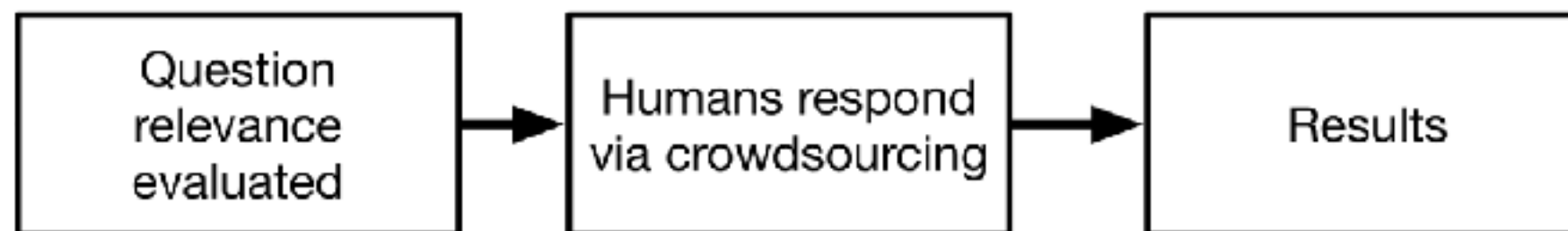
Data collection done with crowdsourcing can be expensive



Goal: achieve goal with minimal responses

One solution

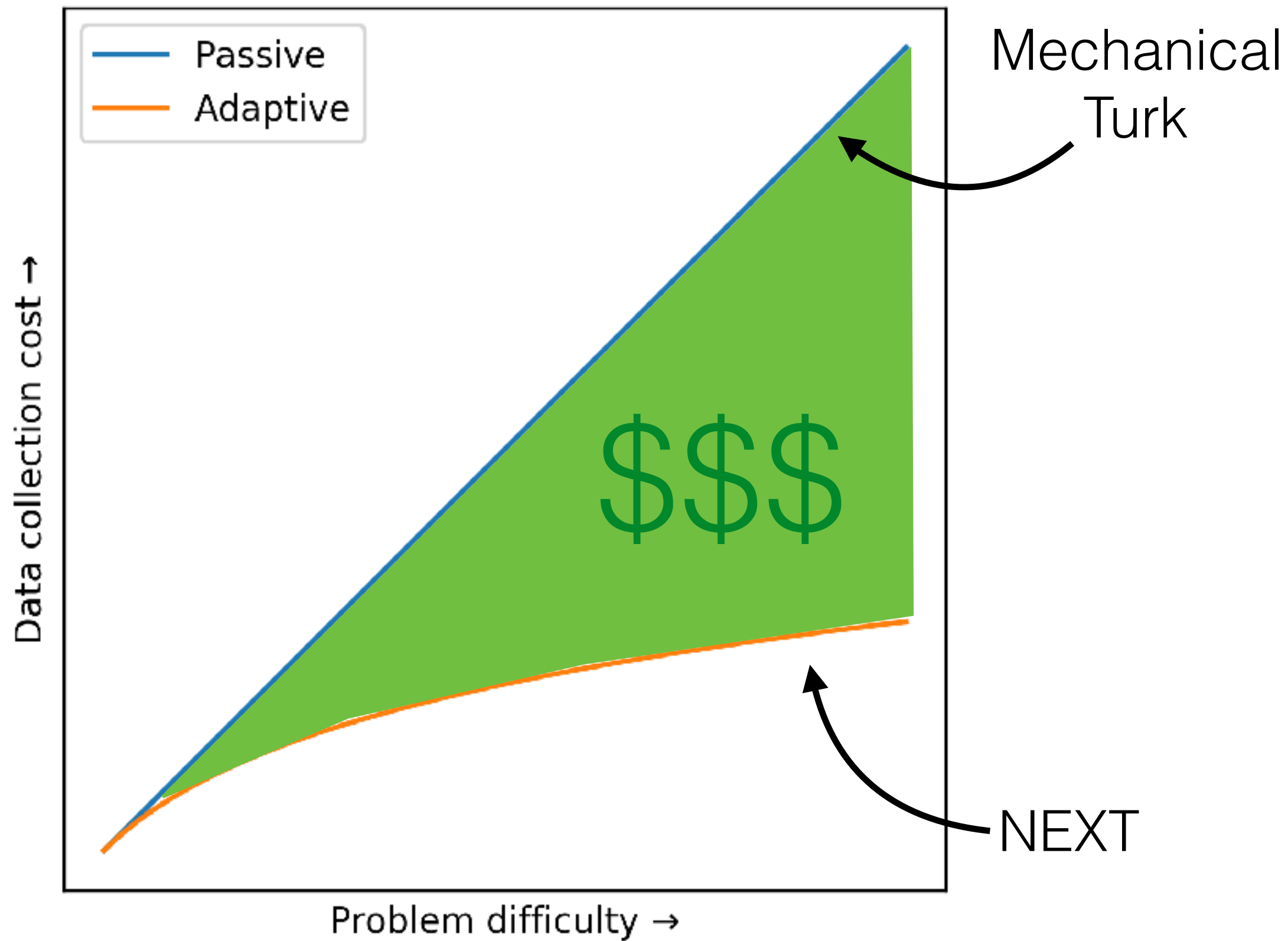
Existing crowdsourcing systems are *passive*



Adapting to previous responses requires fewer data

Goal: adapt to previously collected responses

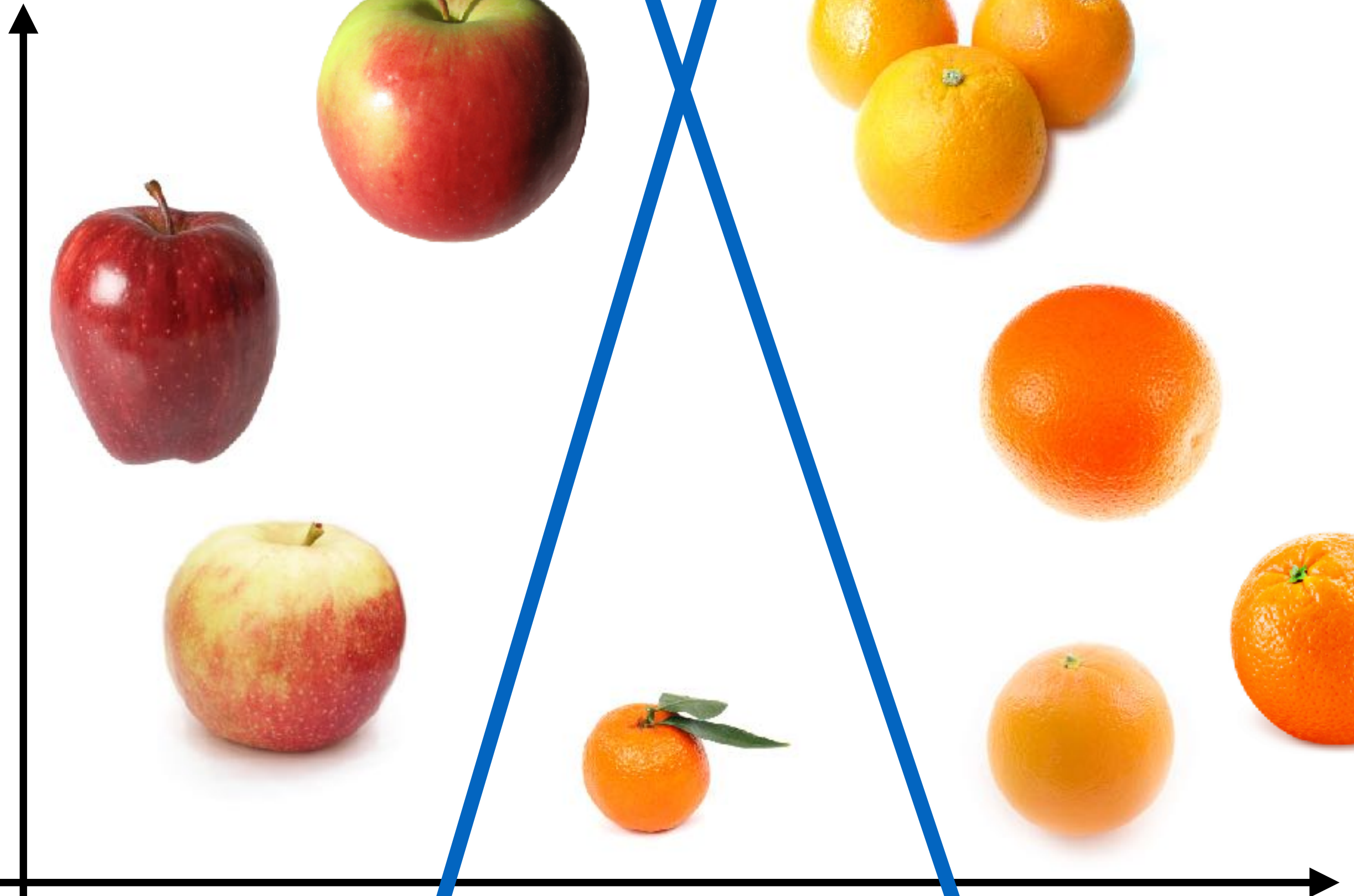
Benefits

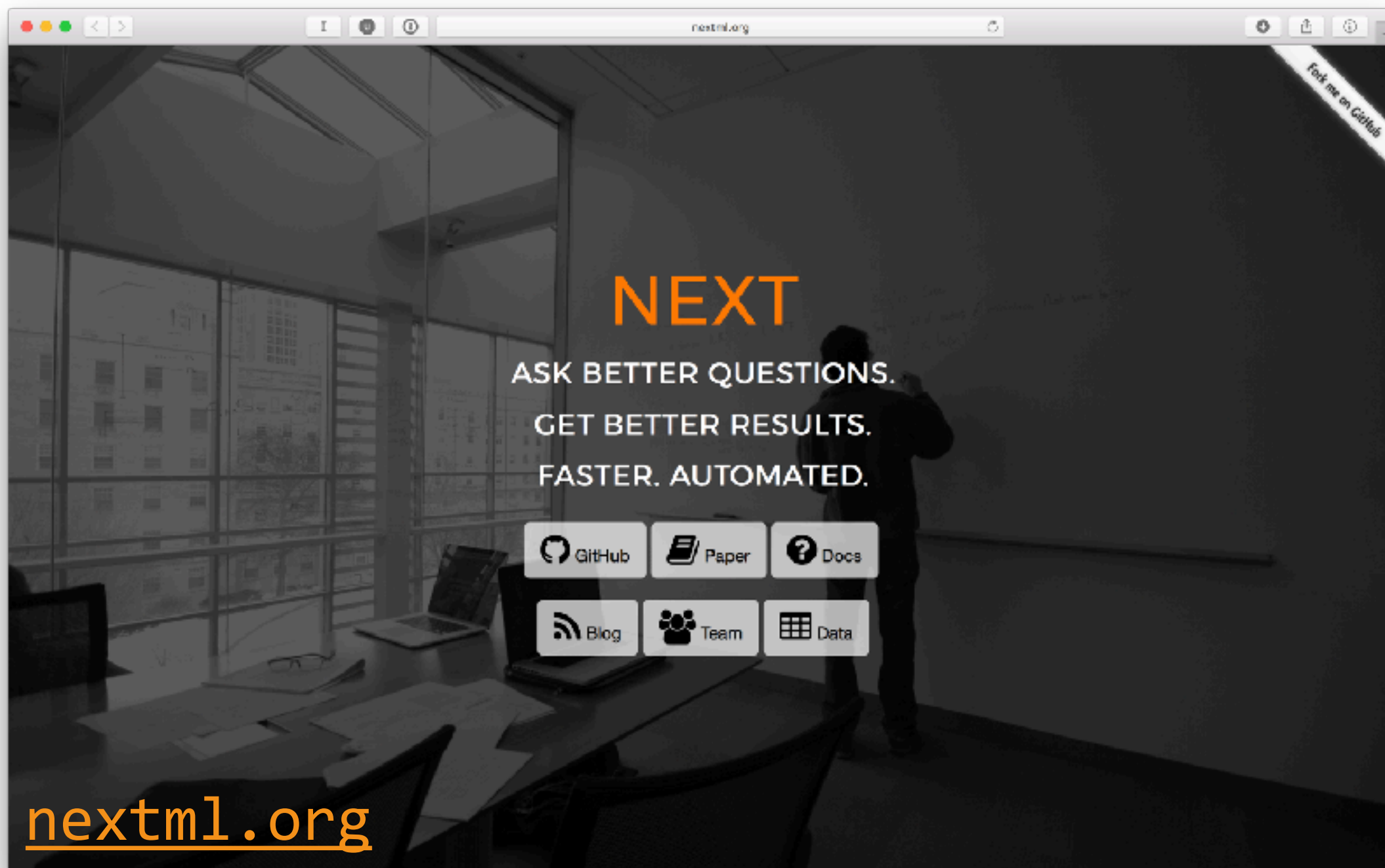


Adaptive sampling can have large benefits

Example solution

Adapting to previous responses yields better results





[Lalit Jain](#)



Daniel Ross



[Rob Nowak](#)



[Kevin Jamieson](#)

Homepage: <http://nextml.org>

Source: <https://github.com/nextml/NEXT>

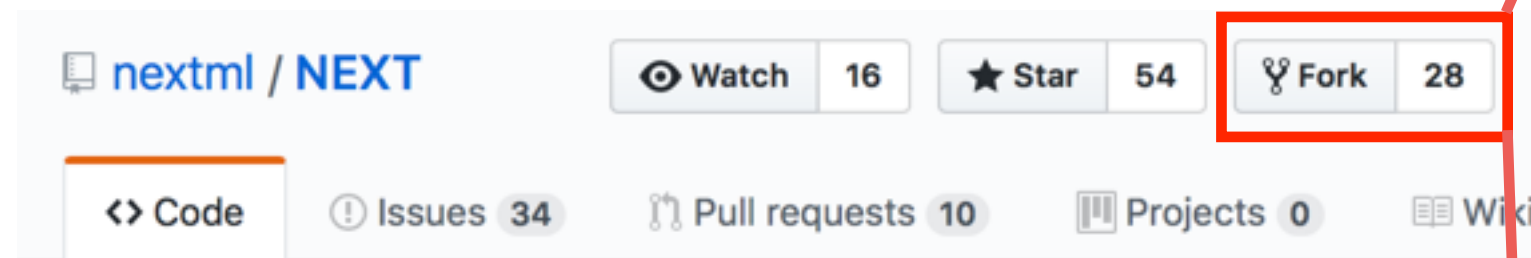
Documentation: <https://github.com/nextml/NEXT/wiki>



tinyurl.com/scipy-next

NEXT users

Theory

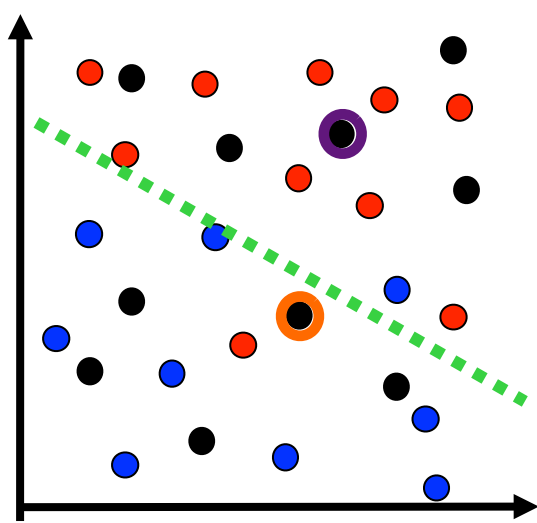


Practice

- nextml / NEXT
- aashish24 / NEXT
- abiswas3 / NEXT
- alphaprime / NEXT
- aniruddhajb / NEXT
- AvinWangZH / NEXT
- ayonsn017 / NEXT
- caomw / NEXT
- connectthefuture / NEXT
- crcox / NEXT
- dconathan / NEXT
- robinsonkwa / NEXT
- jattenberg / NEXT
- jimwmg / NEXT
- justicelee / NEXT
- juthawong / NEXT
- liamim / NEXT
- mllewis / NEXT
- NandanaSengup / NEXT
- naveendennis / NEXT
- pedmiston / NEXT
- samim23 / NEXT
- stsievert / NEXT
- BhargavaA / NEXT
- suchow / NEXT
- sumeetsk / NEXT
- widoptimization / NEXT
- worldbank / NEXT

ML Researchers

Air Force Research Lab uses NEXT for active image classification.



Experimentalists

UW Psychology uses NEXT to find the best algorithms for adaptive data collection in cognitive science.



Practitioners

The New Yorker uses NEXT to crowd-source the weekly cartoon caption contest.



Example problem

THE NEW YORKER



YOUR CAPTION

Enter your caption (250 characters or fewer):

The New Yorker has to find the funniest caption from ~5,000 captions

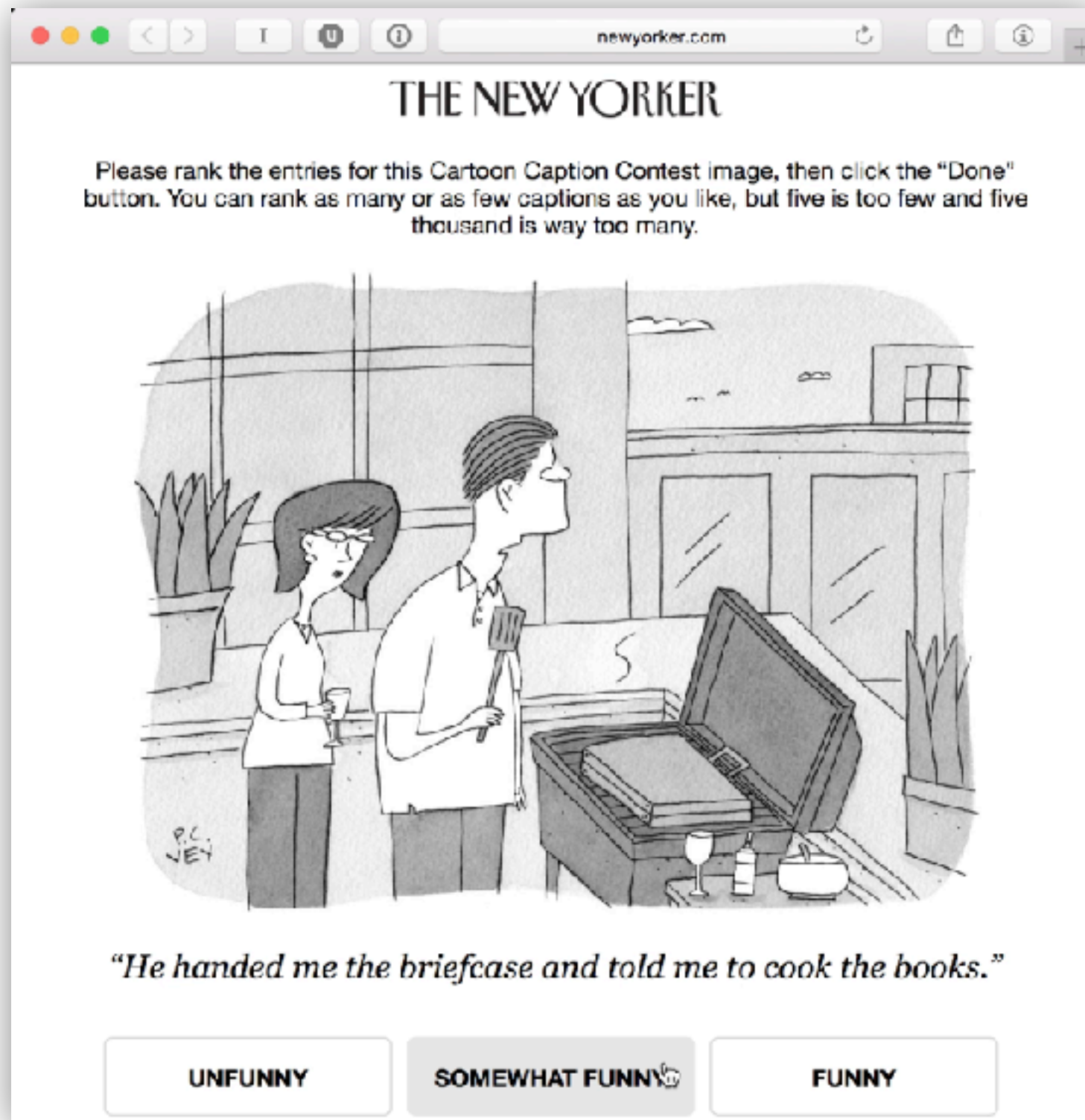


[Bob Mankoff](#)

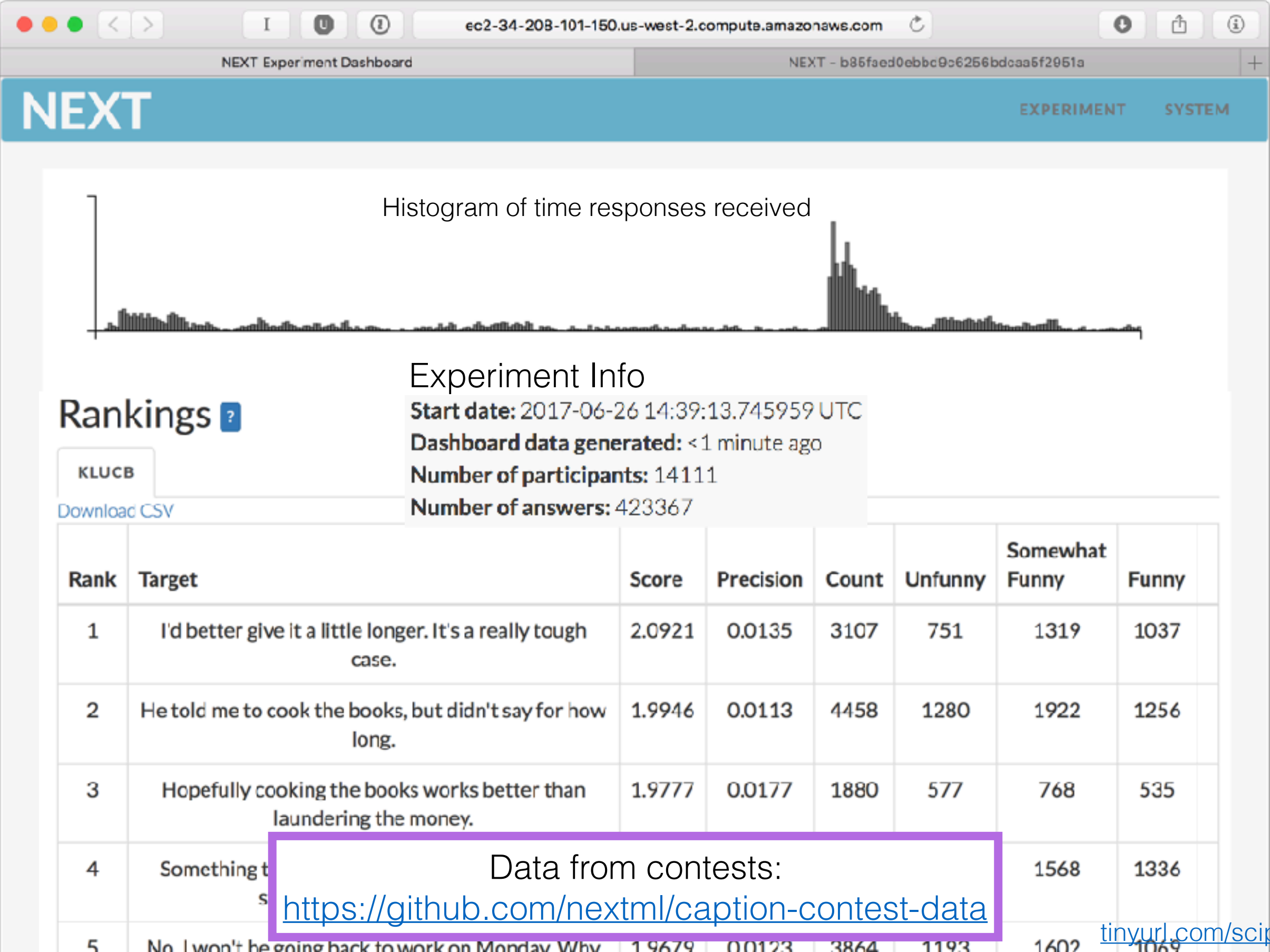
Interface

<http://www.newyorker.com/cartoons/vote>

<http://nextml.org/captioncontest>



Dashboard



Goal

Adaptive
sampling
algorithms

Crowdsourcing

fewer responses
more accurate models

real-world data
participant fatigue
algorithm delays
participant label quality

Goal: let both parties *easily* use NEXT

Software uses

By default, NEXT can be applied to 3 problems

Select face on the bottom most similar to the face on top



Pool based triplets

Cardinal Bandits



comic by [P. C. Vey](#)

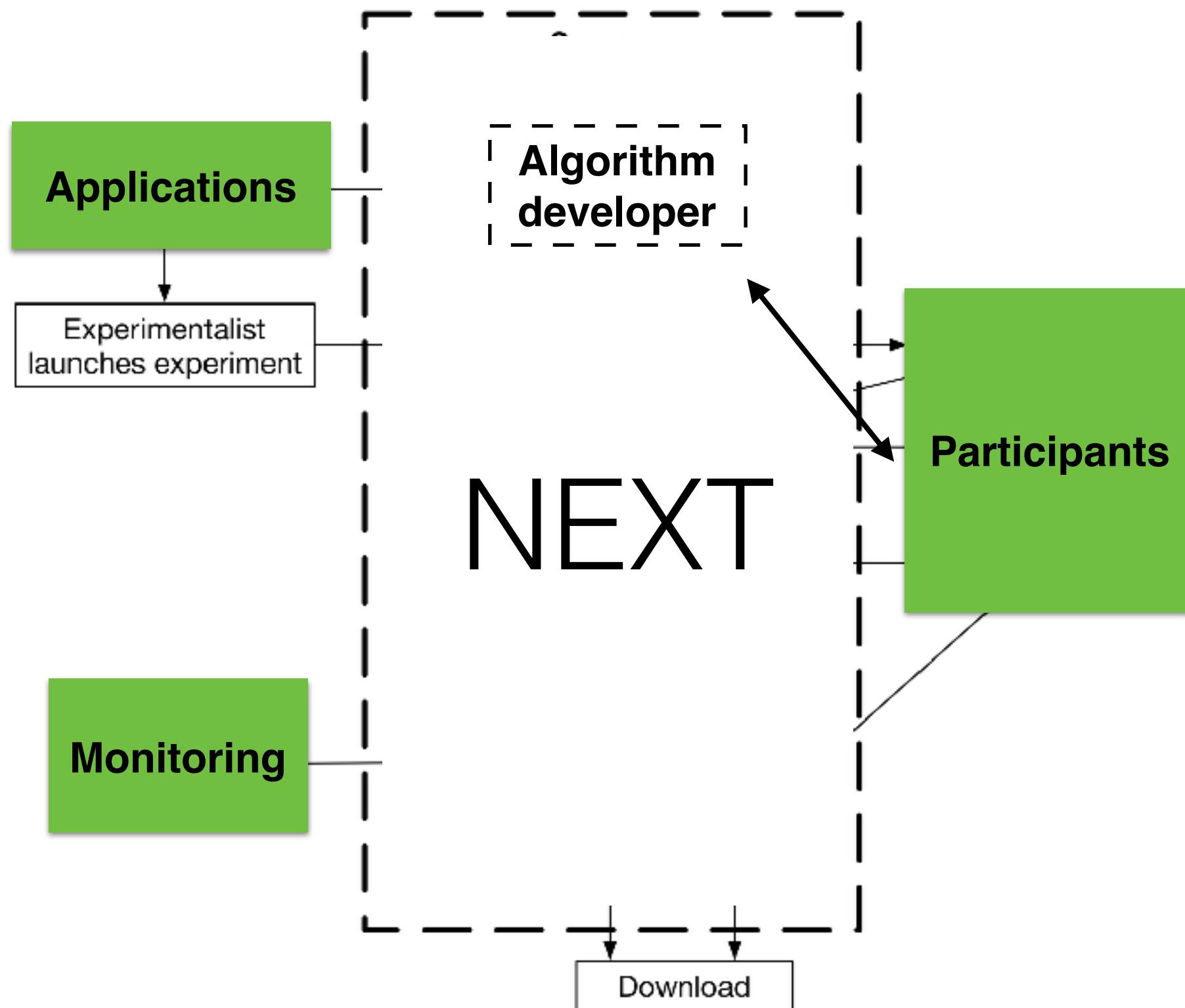
Dueling Bandits

Select the street that looks safer

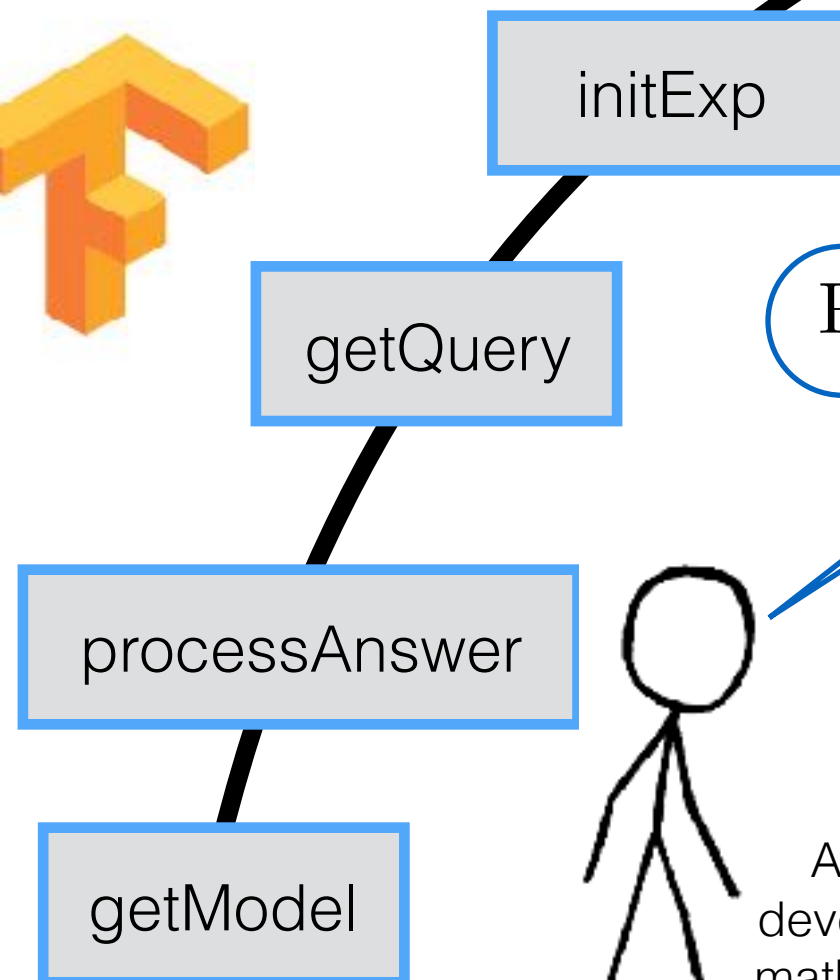
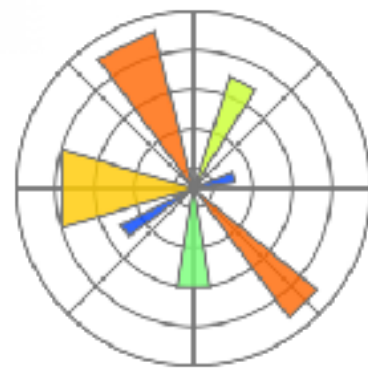
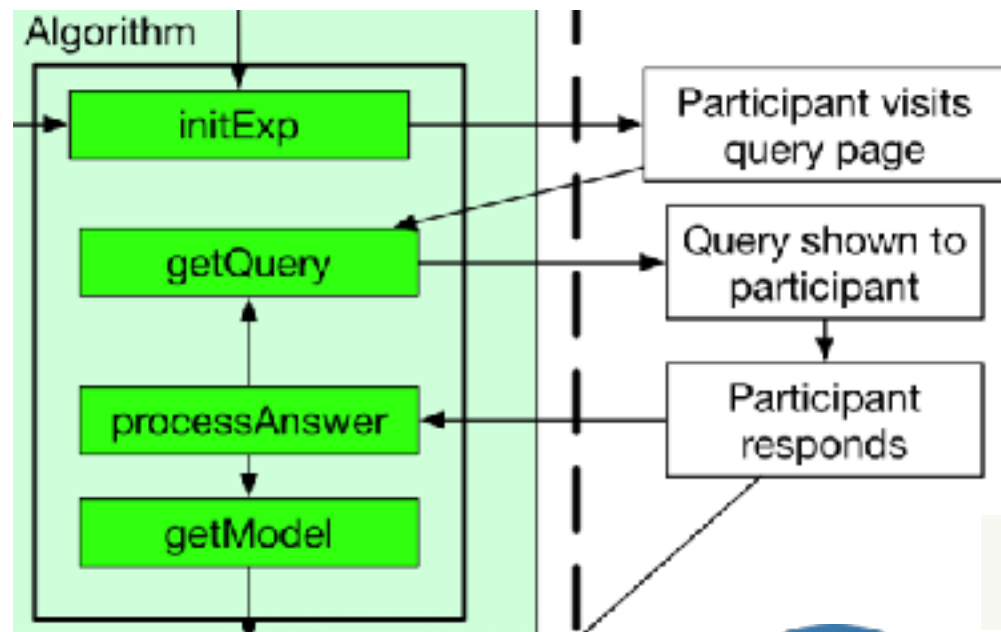


NEXT can also be used with REST API

NEXT use



Algorithm developer use



$$\Pr (|y - \hat{y}| < \epsilon) \leq 1 - \delta$$



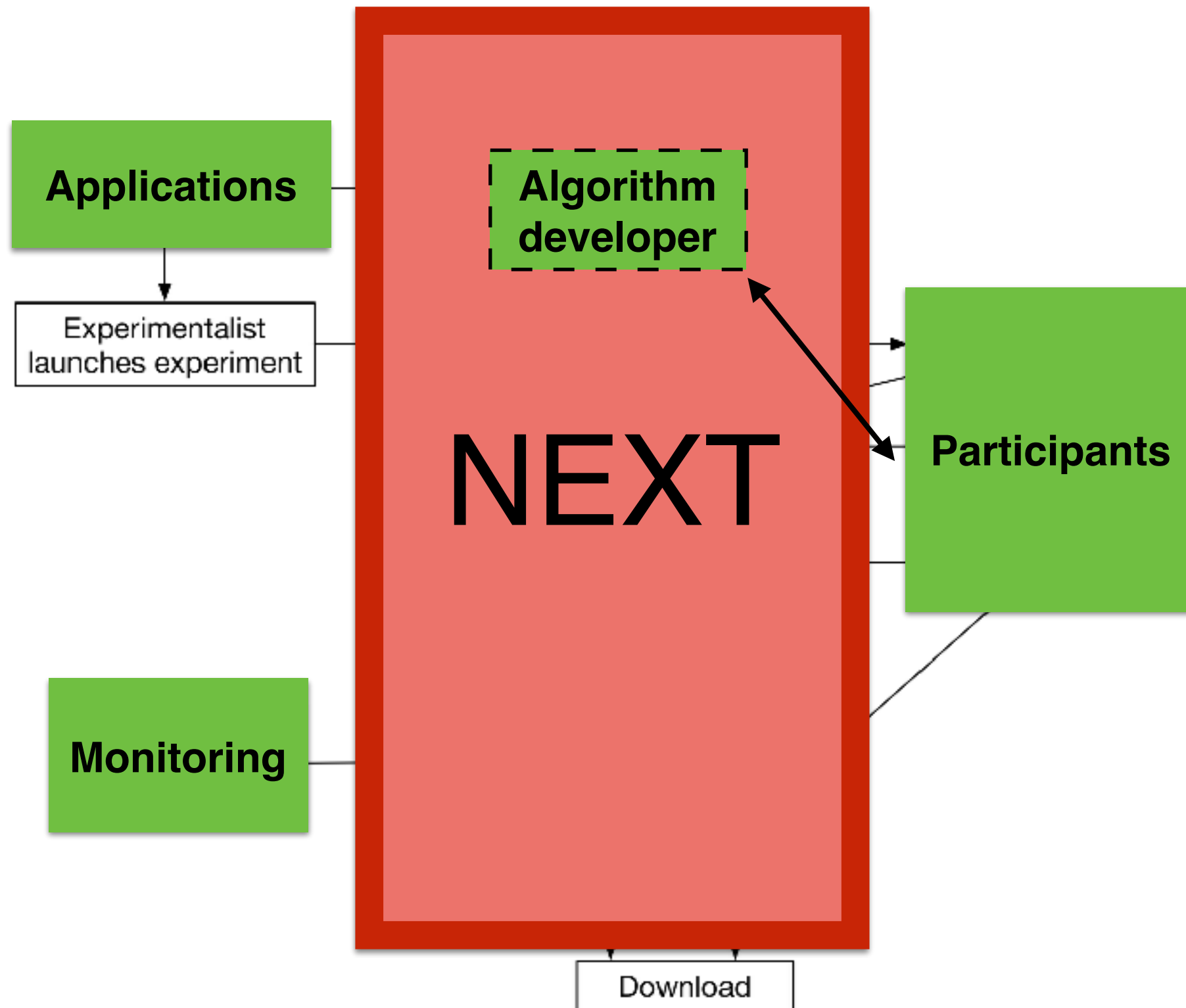
Algorithm
developer and
mathematician

Algorithm design decisions

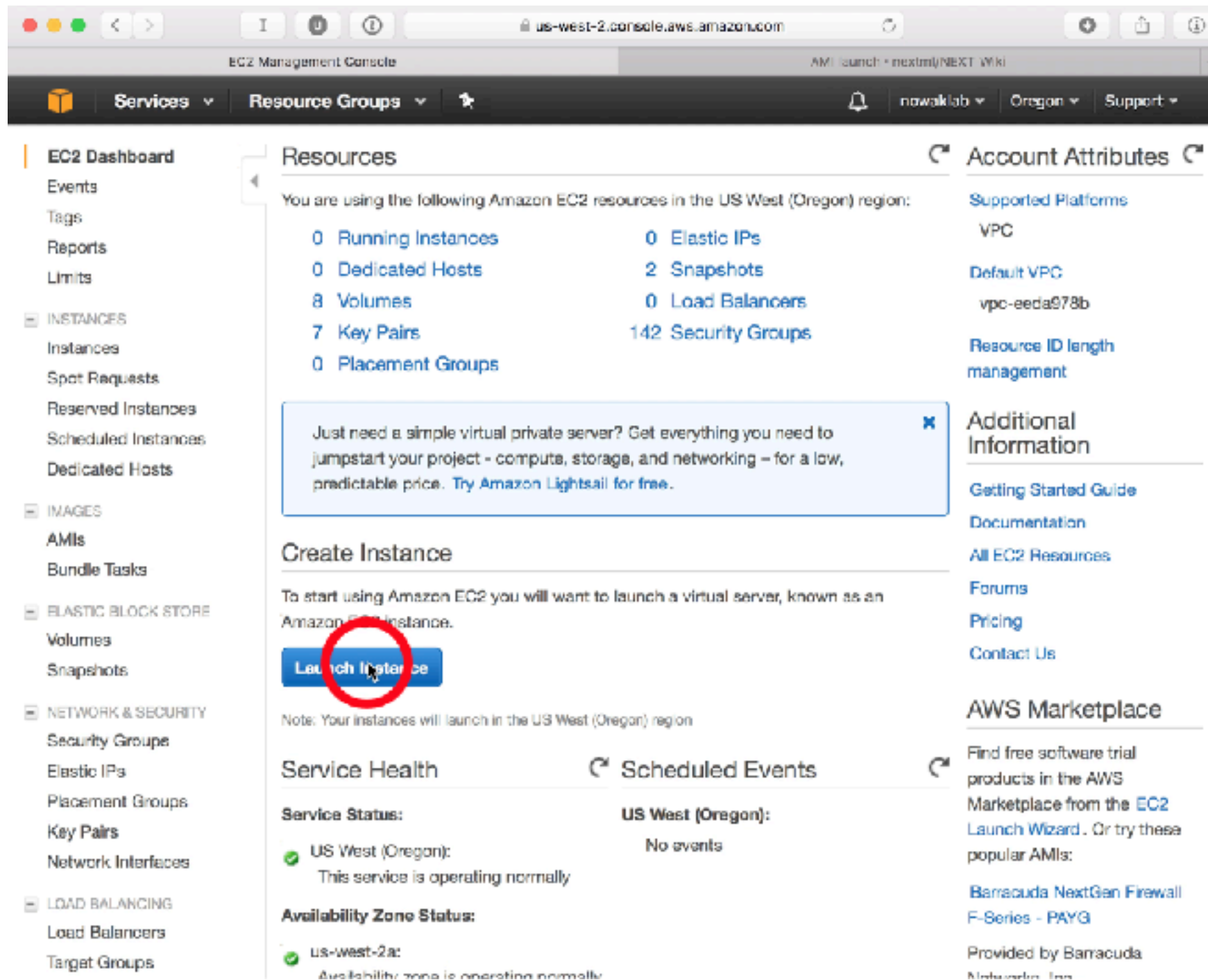
1. Treat algorithms as black boxes
 - (for each function, inputs and outputs are documented and type-checked)
2. Use wrapper to allow *easy* access to experiment information and background jobs
3. Objects are abstracted to integers (i.e., object 42, not `{ 'filename': 'foo.png', 'url': ... }`)

(more detail in proceedings and on docs)

Algorithm use



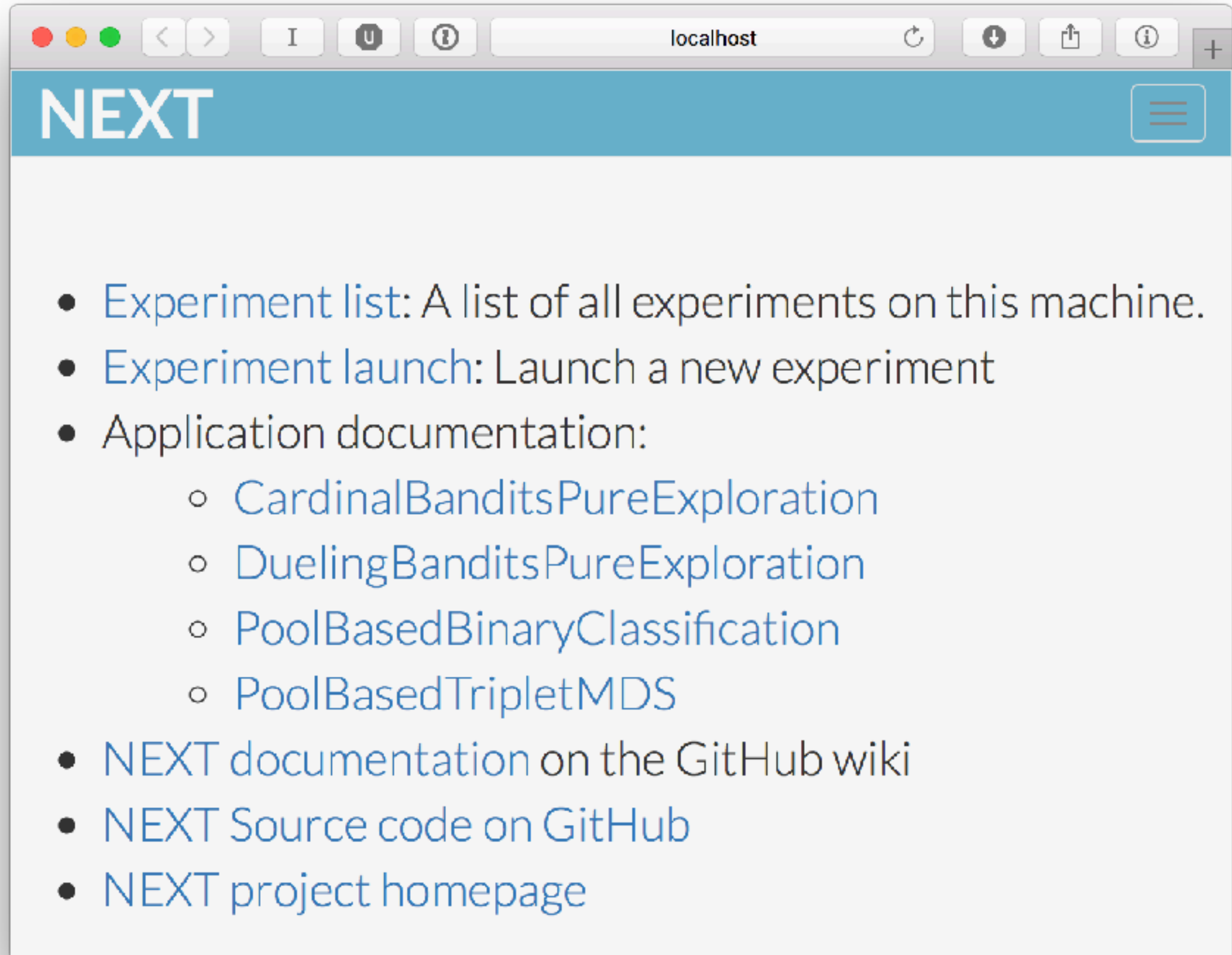
Launching NEXT via Amazon EC2 AMI



(more detail in proceedings
and on docs)

See <https://github.com/nextml/NEXT/wiki>
for details and more launching options

NEXT startup page



Key messages

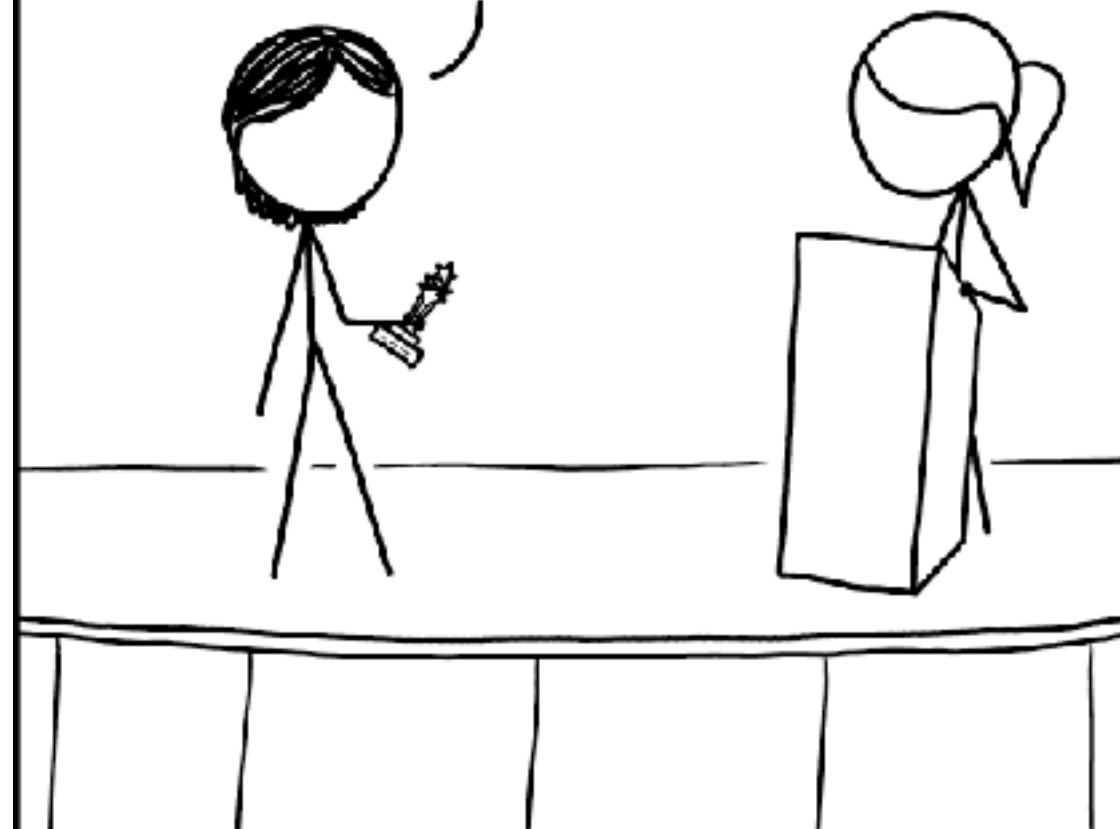
1. Adaptive sampling reduces data collection cost.
2. NEXT is a crowdsourcing data collection tool that can use adaptive sampling techniques
3. NEXT is easy* to use by experimentalists, algorithm developers and practitioners, and a mathematical background is not required.
4. NEXT developers experimentalist engagement to aid research and to gain feedback to improve the software

* NEXT has been created by an academic research group for collaboration with psychologists

I'D LIKE TO THANK MY DIRECTOR,
MY FRIENDS AND FAMILY, AND—
OF COURSE—THE WRITHING MASS
OF GUT BACTERIA INSIDE ME.

I MEAN, THERE'S LIKE ONE OR
TWO PINTS OF THEM IN HERE;
THEIR CELLS OUTNUMBER MINE!

ANYWAY, THIS WAS A
REAL TEAM EFFORT.



Extras...

Algorithm inputs and outputs

- Documented exactly in `apps/[app-id]/algs/Algs.yaml`

```
getQuery:
  args:
    participant_uid:
      type: string
      description: ID of the participant answering the query
  rets:
    description: The index of the target to ask about
    type: num
```

- Function implementation

```
import random

def getQuery(self, butler, participant_uid):
    n = butler.algorithms.get(key='n')
    return random.choice(n)
```

Depends on a library we developed:

<https://github.com/daniel3735928559/pijemont>

Adaptive data flow

- Fundamentally requires 4 functions:
 - initExp**: experiment initialization
 - getQuery**: selects query to present to participant
 - processAnswer**: process participants response
 - getModel**: provides experiment monitoring

