

# CONTEXTUAL BANDIT ALGORITHMS FOR INTERNET-SCALE APPLICATIONS



**LEV REYZIN**

**UIC MATHEMATICS**

**ASA: RECENT ADVANCES IN MACHINE LEARNING**

- Mail
- News
- Finance
- Sports
- Movies
- omg!
- Shine
- Autos
- Shopping
- Travel
- Dating
- Jobs

More Y! Sites >

Make YAHOO! your homepage



Netflix-Try for free. Instantly watch TV shows and movies on your TV or computer. Get 1 month free.



### Rare Yosemite spectacle caught on film

A photographer snaps a once-in-a-lifetime image of a bolt of lightning cracking through a rainbow. [See this picture »](#)

6 - 10 of 80

--	--	--	--	--

All Stories News Local Entertainment Sports More >



### Rod Stewart Plays Intimate L.A. Club Gig for Diehard Fans, Harry Styles

On May 7, rock icon Rod Stewart, who experienced a career rebirth in the 1990s and 2000s with his multi-volume Great American Songbook Maximum Performance (NEW)



### Beyonce and Blue Ivy Step Out in Paris, Gwyneth Paltrow Wears a Sheer Dress and No Underwear: Today's Top Stories

Take a look at Us Weekly's most-read stories from Thursday, April 25 Us Weekly

Cycling-UCI go on the attack after latest accusations by USADA

### Trending Now

Watch the show >

- |                           |                          |
|---------------------------|--------------------------|
| 1 Jessica Alba double...  | 6 Psy knocked off        |
| 2 Ace of Base Nazi past   | 7 Danica Patrick divo... |
| 3 38 die in mental hos... | 8 T-Mobile 'deceptive'   |
| 4 Jennifer Love Hewitt    | 9 HIV vaccine fails      |
| 5 Edward Norton fight     | 10 Hyundai suicide ad    |

Ad Feedback

AdChoices

Mount Laurel

48°F Fair



Today 65° 40°



Tomorrow 60° 44°



Sunday 60° 44°

- Mail
- News
- Finance
- Sports
- Movies
- omg!
- Shine
- Autos
- Shopping
- Travel
- Dating
- Jobs

More Y! Sites >

Make YAHOO! your homepage



Netflix-Try for free. Instantly watch TV shows and movies on your TV or computer. Get 1 month free.



### Rare Yosemite spectacle caught on film

A photographer snaps a once-in-a-lifetime image of a bolt of lightning cracking through a rainbow. [See this picture »](#)

6 - 10 of 80

--	--	--	--	--

All Stories News Local Entertainment Sports More >

**Rod Stewart Plays Intimate L.A. Club Gig for Diehard Fans, Harry Styles**  
 On May 7, rock icon Rod Stewart, who experienced a career rebirth in the 1990s and 2000s with his multi-volume Great American Songbook Maximum Performance (NEW)

**Beyonce and Blue Ivy Step Out in Paris, Gwyneth Paltrow Wears a Sheer Dress and No Underwear: Today's Top Stories**  
 Take a look at Us Weekly's most-read stories from Thursday, April 25  
 Us Weekly

**Cycling-UCI go on the attack after latest accusations by USADA**

Trending Now [Watch the show »](#)

1 Jessica Alba double...	6 Psy knocked off
2 Ace of Base Nazi past	7 Danica Patrick divo...
3 38 die in mental hos...	8 T-Mobile 'deceptive'
4 Jennifer Love Hewitt	9 HIV vaccine fails
5 Edward Norton fight	10 Hyundai suicide ad

**EVONY**  
 FREE FOREVER  
 PLAY NOW

Ad Feedback AdChoices

Mount Laurel  
 48°F Fair

Today 65° 40°	Tomorrow 60° 44°	Sunday 60° 54°
------------------	---------------------	-------------------

Watch the Internet of Everything in action.



LATEST HEADLINES

7:48 AM Shiseido Sees Profit Next Year

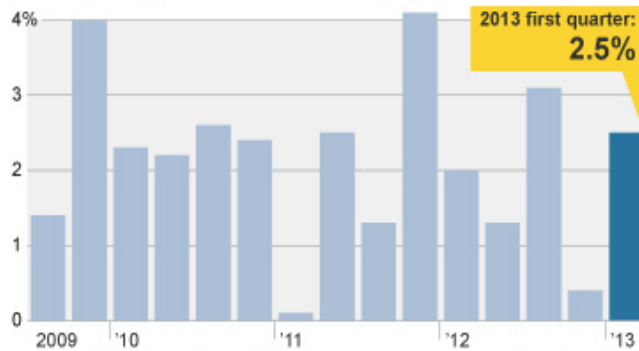
MORE HEADLINES »

## First-Quarter Growth, at 2.5%, Misses Expectations

The nation's gross domestic product, a measure of all goods and services produced in the economy, advanced at a 2.5% annual rate between January and March, the Commerce Department said. [17 min ago](#)

### U.S. GDP

Quarterly change at an annualized rate, adjusted for inflation



### Disappointing GDP Data Weigh on Futures

Mansion: From Oil, a Land Rush

Einstein Proved Right--Again

NFL Draft

### Markets >

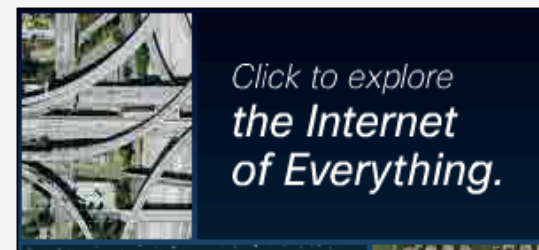
Sponsored by: Vanguard®

Overview U.S. Europe Asia FX Rates Futures

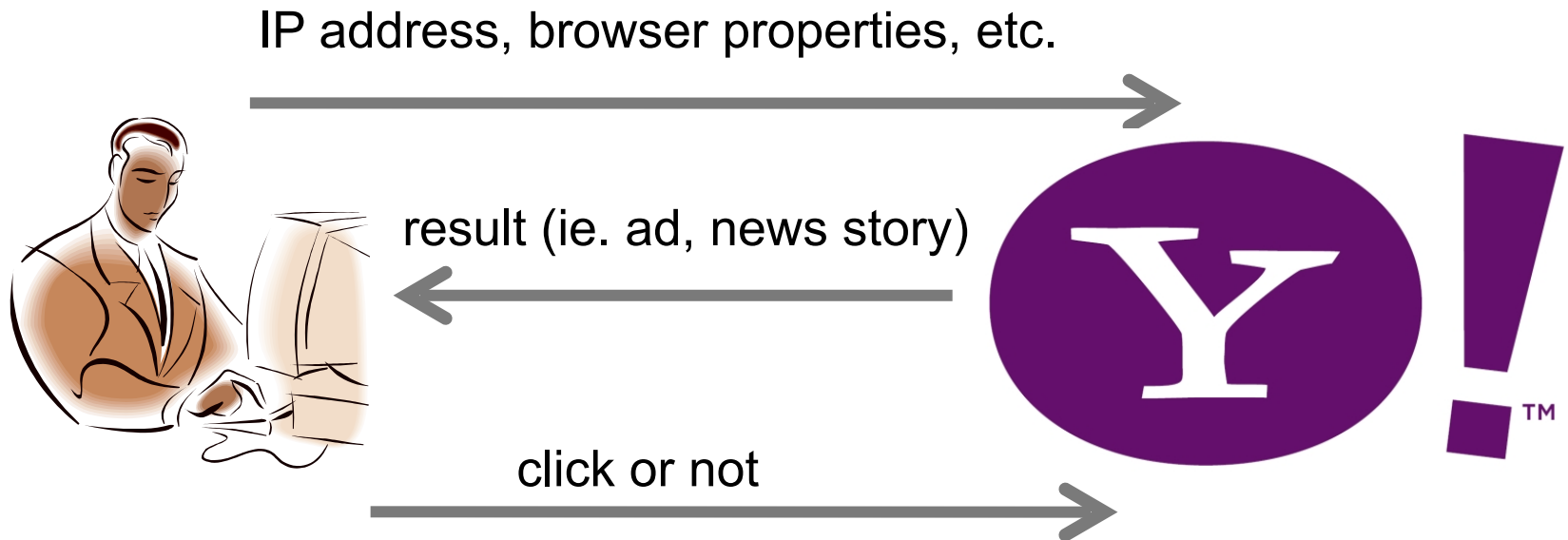
	LAST	CHG	%CHG	RANGE: 1 DAY
DJIA	14700.80	+24.50	0.17%	1473
Nasdaq	3289.99	+20.33	0.62%	1471
FTSE 100	6401.63	-40.96	0.64%	
Nikkei 225	13884.13	-41.95	0.30%	1463
Crude Oil	93.50	-0.14	0.15%	
Gold	1470.50	+8.50	0.58%	

4/26/13 8:44 AM EDT

Market Data | MoneyBeat | Watchlist | Portfolio | Customize



# SERVING CONTENT TO USERS



# SERVING CONTENT TO USERS



IP address, browser properties, etc.



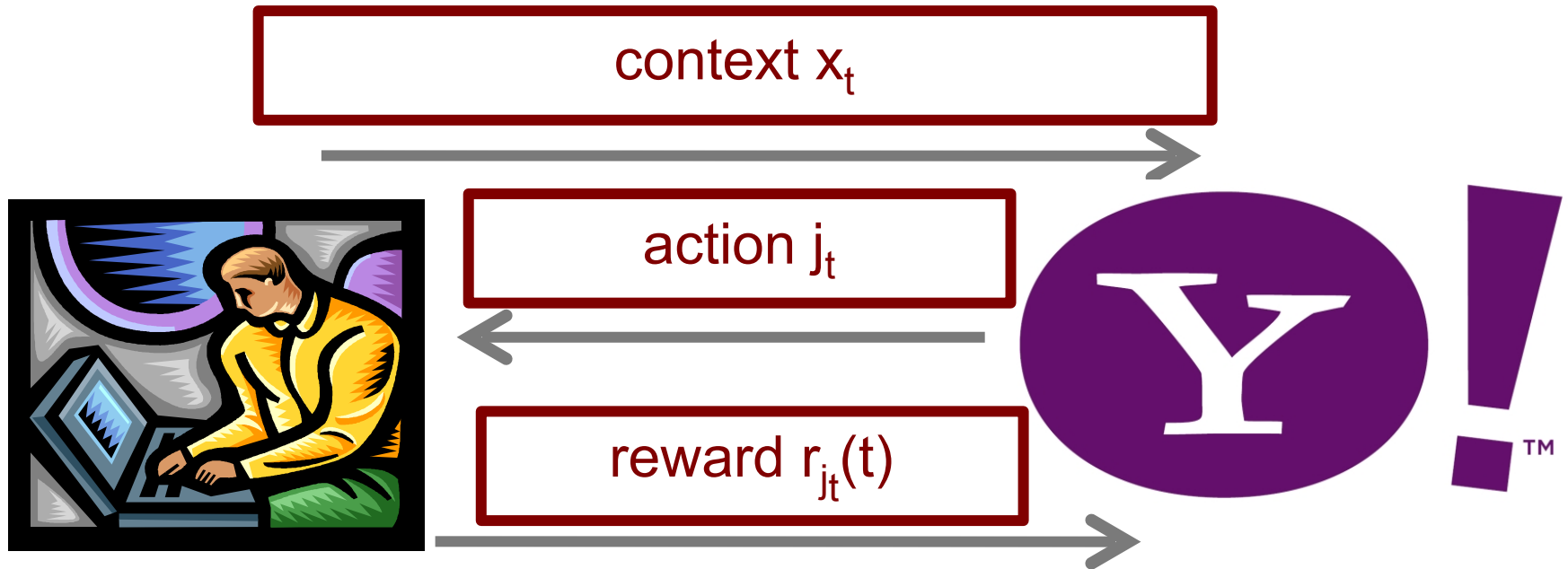
result (ie. ad, news story)



click or not

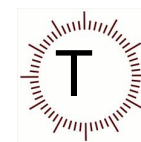
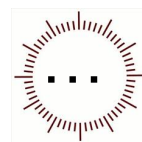


# SERVING CONTENT TO USERS



# MULTIARMED BANDITS

## [ROBBINS '52]



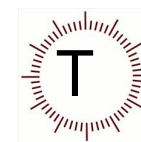
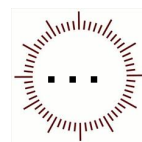
⋮





# MULTIARMED BANDITS

## [ROBBINS '52]



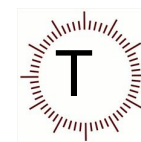
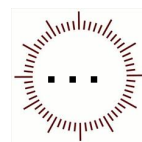
⋮



click

# MULTIARMED BANDITS

## [ROBBINS '52]



click



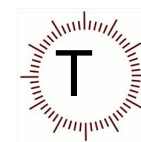
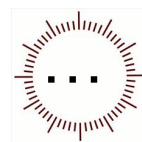
click

⋮



# MULTIARMED BANDITS

## [ROBBINS '52]



⋮



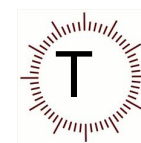
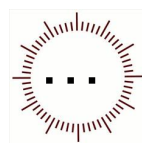
click

no

click

# MULTIARMED BANDITS

## [ROBBINS '52]



click



click

no



click

⋮

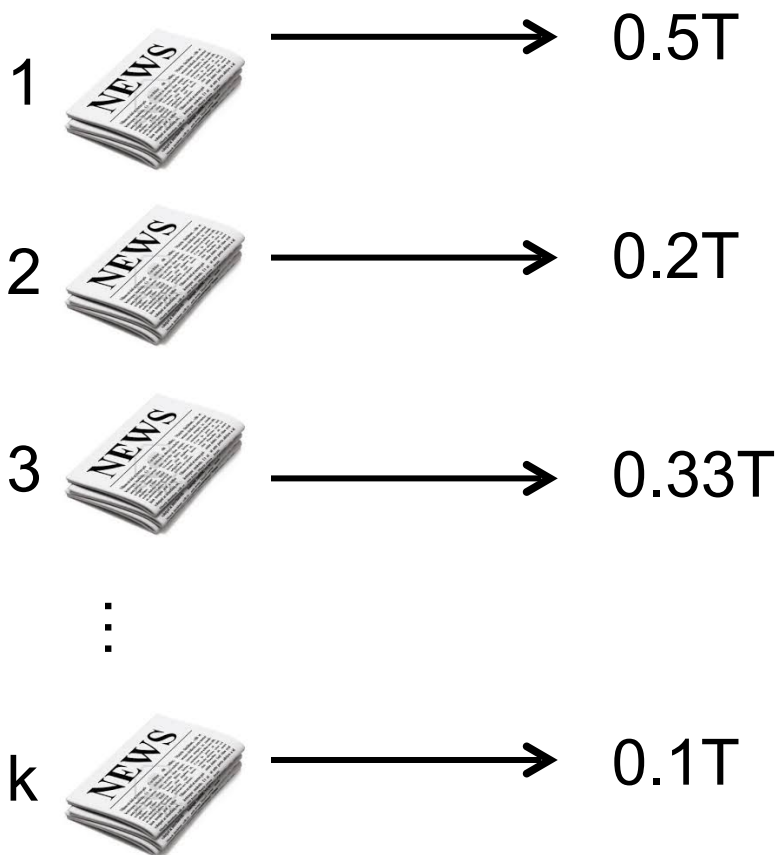


click



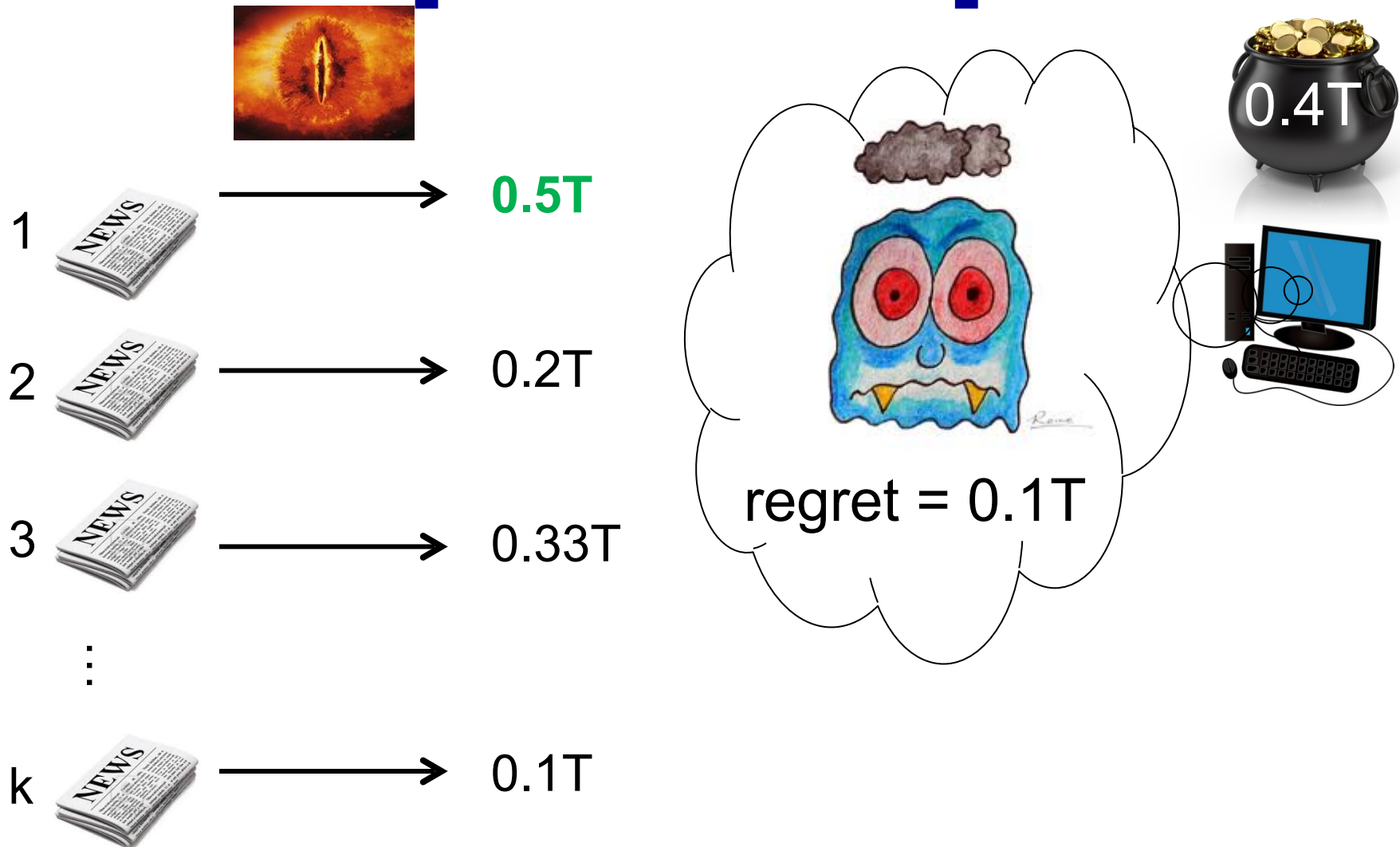
# MULTIARMED BANDITS

## [ROBBINS '52]



# MULTIARMED BANDITS

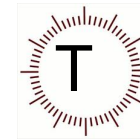
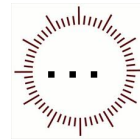
## [ROBBINS '52]



# CONTEXTUAL BANDITS

## [AUER-CESABIANCHI-FREUND-SCHAPIRE '02]

context:



1



2



3



⋮

k

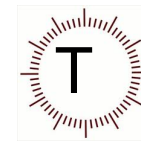
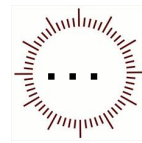


N experts/policies/functions  
think of  $N \gg K$

# CONTEXTUAL BANDITS

## [AUER-CESABIANCHI-FREUND-SCHAPIRE '02]

context:  $x_1$



1



2



3



⋮

k



5



1



1



4



K



3

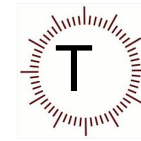
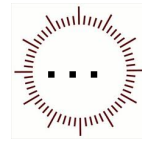
N experts/policies/functions  
think of  $N \gg K$



# CONTEXTUAL BANDITS

## [AUER-CESABIANCHI-FREUND-SCHAPIRE '02]

context:  $x_1$



1



click

2



3



⋮

k



5



1



1



4



K

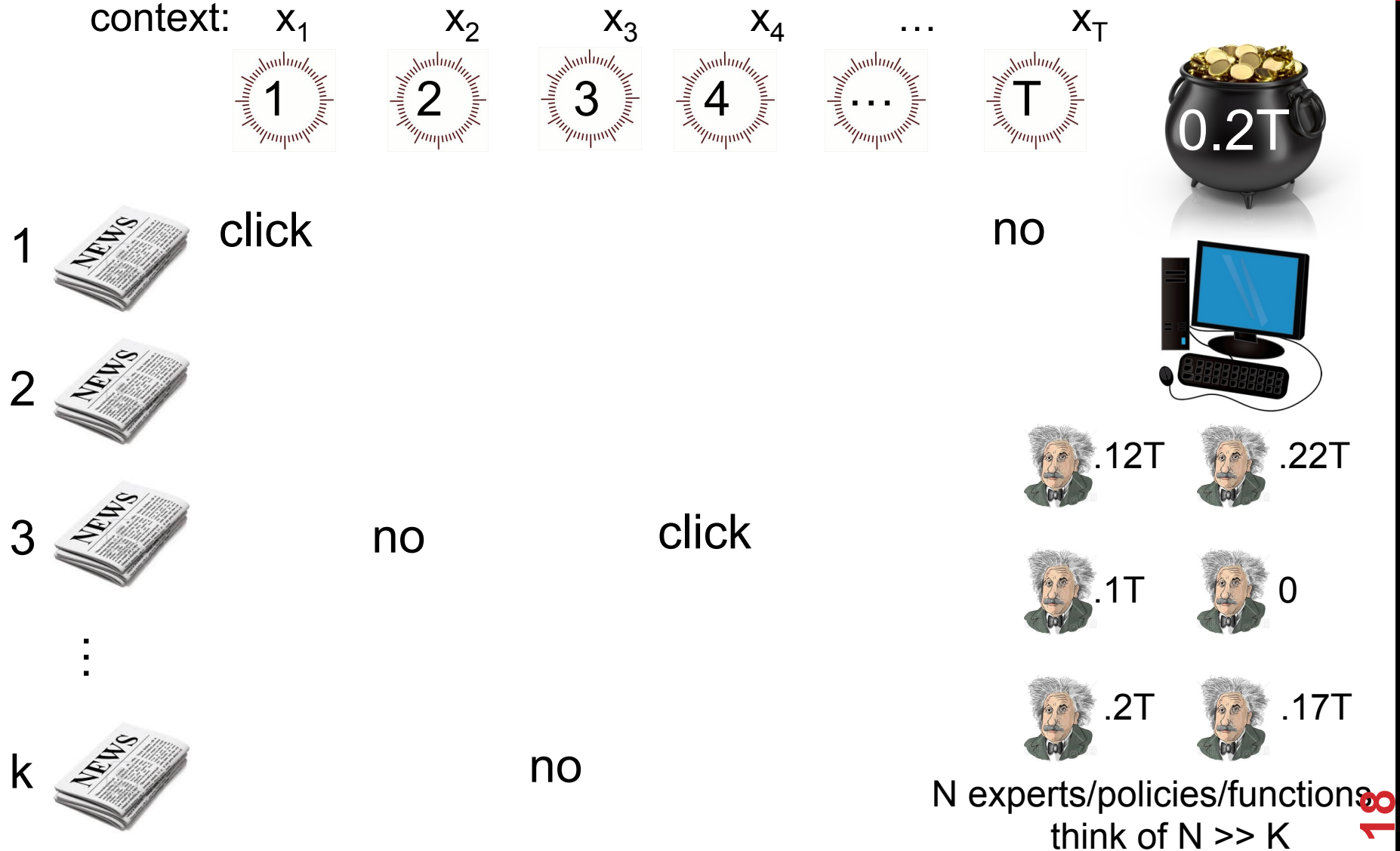


3

N experts/policies/functions  
think of  $N \gg K$

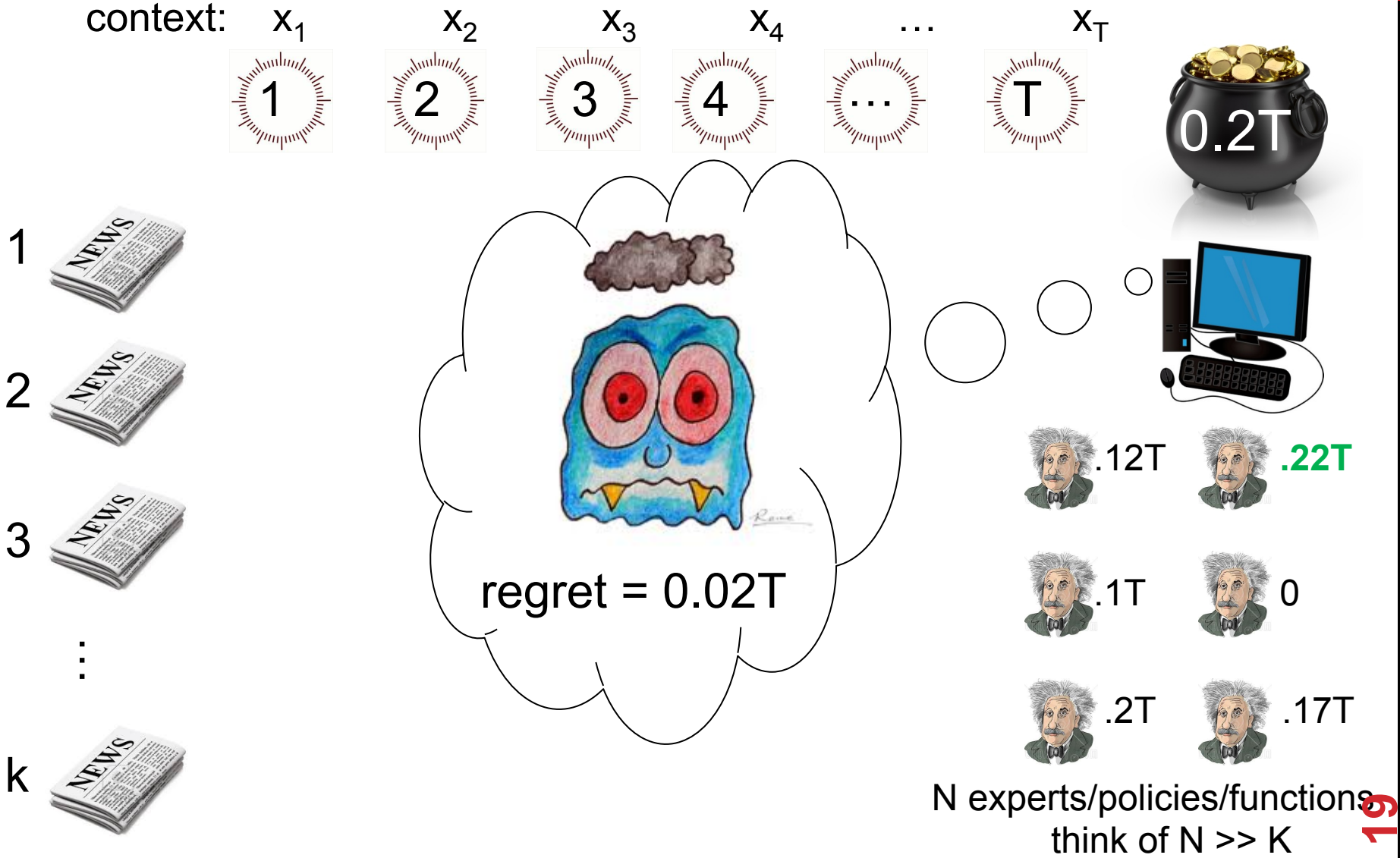
# CONTEXTUAL BANDITS

## [AUER-CESABIANCHI-FREUND-SCHAPIRE '02]



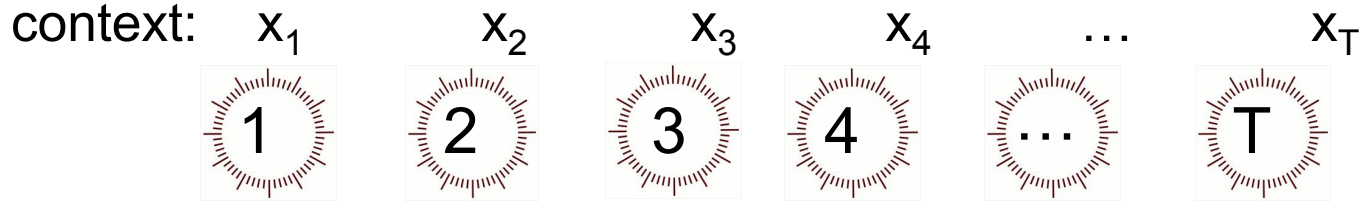
# CONTEXTUAL BANDITS

## [AUER-CESABIANCHI-FREUND-SCHAPIRE '02]



# CONTEXTUAL BANDITS

## [AUER-CESABIANCHI-FREUND-SCHAPIRE '02]



⋮



the clicks can come i.i.d. from a distribution or be arbitrary

stochastic / adversarial

The experts can be present or not.

contextual / non-contextual



# BANDITS

Harder than supervised learning:

In the bandit setting we do not know the rewards of actions not taken.

Many applications

Ad auctions, medicine, finance, ...

Exploration/Exploitation

Can **exploit** expert/article you've learned to be good.

Can **explore** expert/article you're not sure about.

# EPSILON-FIRST

Rough idea of  **$\epsilon$ -first** (or  **$\epsilon$ -greedy**): act randomly for  $\epsilon$  rounds, then go with best (arm or expert).

Rough analysis: even for 2 arms, we suffer regret  $\epsilon + (T-\epsilon)/(\epsilon^{1/2})$ .

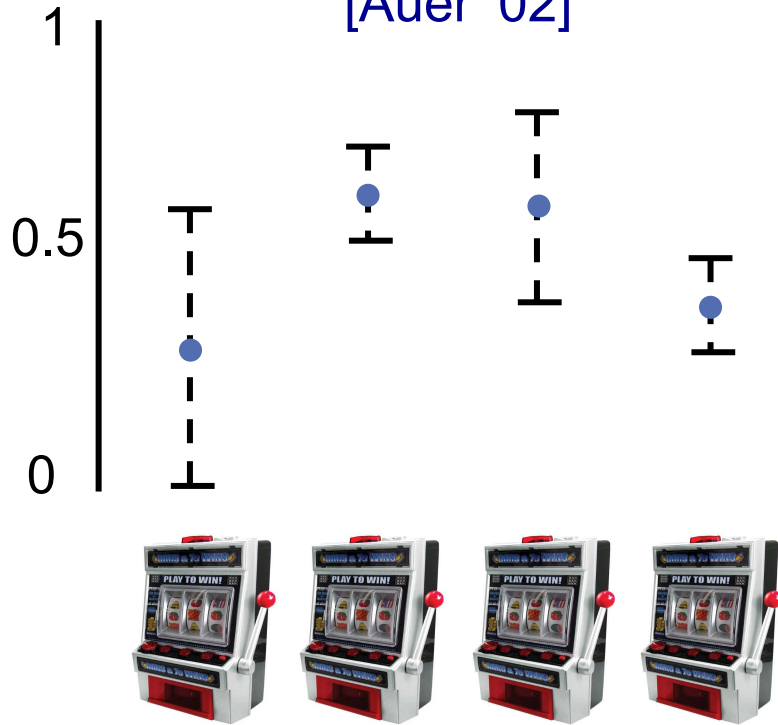
- $\epsilon \approx T^{2/3}$  is optimal tradeoff, gives regret  $\approx T^{2/3}$

**But actually  $O(T)^{1/2}$  regret is possible!**

# TRADITIONAL BANDIT ALGORITHMS

## UCB

[Auer '02]



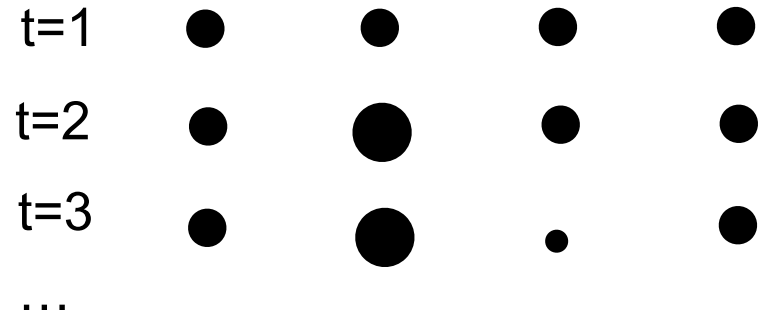
Algorithm: at every time step

- 1) pull arm with highest UCB
- 2) update confidence bound of the arm pulled.

## EXP3 / EW

[Littlestone-Warmuth '94]

[Auer et al. '02]



Algorithm: at every time step

- 1) sample from distribution defined by weights (mixed w/ uniform)
- 2) update weights “exponentially”

# UCB VS EXP3: A COMPARISON

## UCB

[AUER '02]

### ◆ Pros

- ◆ Optimal for the stochastic setting.
- ◆ Succeeds with high probability.

### ◆ Cons

- ◆ Does not work in the adversarial setting.
- ◆ Is not optimal in the contextual setting.

## EXP3 & FRIENDS

[ACFS '02]

### ◆ Pros

- ◆ Optimal for both the adversarial and stochastic settings.
- ◆ Can be made to work in the contextual setting

### ◆ Cons

- ◆ Does not succeed with high probability in the contextual setting (only in expectation).

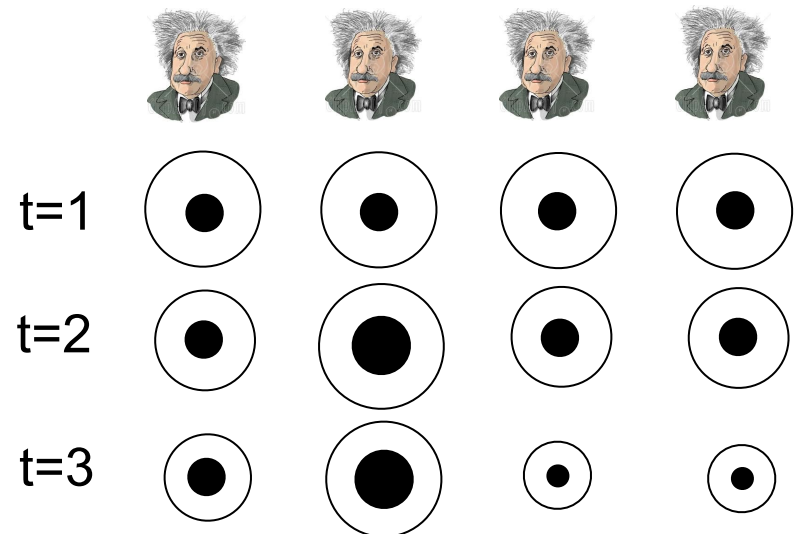


# EXP4.P

**Main Theorem** [Beygelzimer-Langford-Li-R-Schapire '11]:  
For any  $\delta > 0$ , with probability at  $> 1 - \delta$ , EXP4P has “optimal”  
regret in the adversarial contextual bandit setting.

key insights  
on top of UCB/ EXP

- 1) exponential weights and upper confidence bounds “stack”
- 2) generalized Bernstein’s inequality for martingales



# IDEAS BEHIND EXP4.P

(ALL APPEARED IN PREVIOUS ALGORITHMS)

## exponential weights

- keep a weight on each expert that drops exponentially in the expert's (estimated) performance

## upper confidence bounds

- use an upper confidence bound on each expert's estimated reward

## ensuring exploration

- make sure each action is taken with some minimum probability

## importance weighting

- give rare events more importance to keep estimates unbiased

# Exponential Weight Algorithm for Exploration and Exploitation with Experts

## Exp4.P [Beygelzimer, Langford, Li, R, Schapire '10]

Initialization:  $\forall \pi \in \Pi : w_t(\pi) = 1$

For each  $t = 1, 2, \dots$ :

1. Observe  $x_t$  and let for  $a = 1, \dots, K$

$$p_t(a) = (1 - K\rho_{\min}) \frac{\sum_{\pi} \mathbf{1}[\pi(x_t) = a] w_t(\pi)}{\sum_{\pi} w_t(\pi)} + \rho_{\min},$$

where  $\rho_{\min} = \sqrt{\frac{\ln |\Pi|}{KT}}$ .

2. Draw  $a_t$  from  $p_t$ , and observe reward  $r_t(a_t)$ .
3. Update for each  $\pi \in \Pi$

$$w_{t+1}(\pi) = w_t(\pi) \exp \left( \frac{\rho_{\min}}{2} \left( \mathbf{1}[\pi(x_t) = a_t] \frac{r_t(a_t)}{p_t(a_t)} + \frac{1}{p_t(\pi(x_t))} \sqrt{\frac{\ln N/\delta}{KT}} \right) \right)$$

# EXP4.P IN PRACTICE

- ◆ Application – Yahoo! front page
- ◆ We chose a special policy class for which we could efficiently keep track of the weights.
  - ◆ Created 5 clusters, with users (at each time step) getting features based on their distances to clusters.
  - ◆ Policies mapped clusters to article (action) choices.
  - ◆ Ran on personalized news article recommendations for Yahoo! front page.
- ◆ We used a **learning bucket** on which we ran the algorithms and a **deployment bucket** on which we ran the greedy (best) learned policy.

Reported estimated (normalized) click-through rates on front page news. Over **41M user visits**. 253 total articles. 21 candidate articles per visit.

	<b>EXP4P</b>	<b>EXP4</b>	<b><math>\epsilon</math>-greedy</b>
Learning eCTR	1.0525	1.0988	1.3829
Deployment eCTR	<b>1.6512</b>	1.5309	1.4290

Reported estimated (normalized) click-through rates on front page news. Over **41M user visits**. 253 total articles. 21 candidate articles per visit.

	<b>EXP4P</b>	<b>EXP4</b>	<b><math>\epsilon</math>-greedy</b>
Learning eCTR	1.0525	1.0988	1.3829
Deployment eCTR	<b>1.6512</b>	1.5309	1.4290

**Why does this work in practice?**

# HOPE FOR AN EFFICIENT ALGORITHM?

[DUDIK-HSU-KALE-KARAMPATZIAKIS-LANGFORD-R-ZHANG '11]

For EXP4P, the dependence on  $N$  in the regret is logarithmic.

this suggests

We could compete with a large, even super-polynomial number of policies! (e.g.  $N=K^{100}$  becomes  $10 \log^{1/2} K$  in the regret)

however

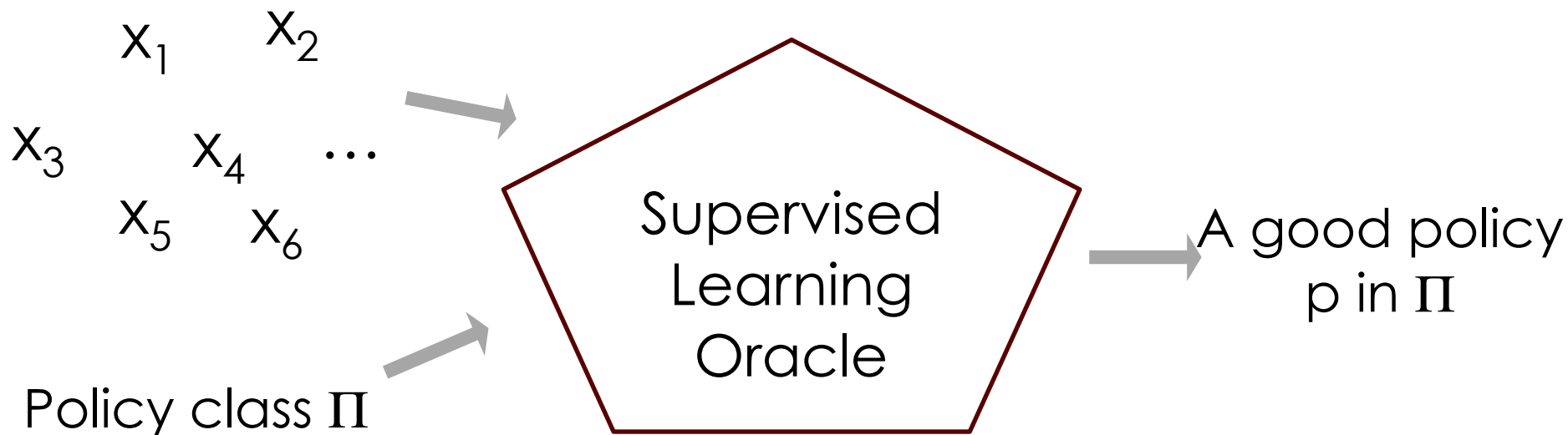
All known contextual bandit algorithms explicitly “keep track” of the  $N$  policies. Even worse, just reading in the  $N$  would take too long for large  $N$ .

# Reduce to Supervised Learning!

32

(Idea from [Langford-Zhang '07])

- ◆ “Competing” with an exponentially large set of policies is **commonplace** in **supervised learning**.
- ◆ Recommendations of the policies/functions don't need to be explicitly read when the policy class has **structure**!





# Reduce to Supervised Learning!

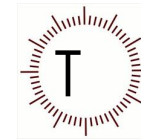
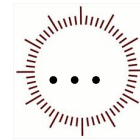
33

(Idea from [Langford-Zhang '07])

- ◆ “Competing” with an exponentially large set of policies is **commonplace** in **supervised learning**.
- ◆ Recommendations of the policies/functions don't need to be explicitly read when the policy class has **structure**!



context:

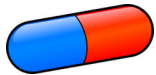
 $x_1$  $x_2$  $x_3$ 

1



yes

2



no

3



⋮

k



5



1



1



4

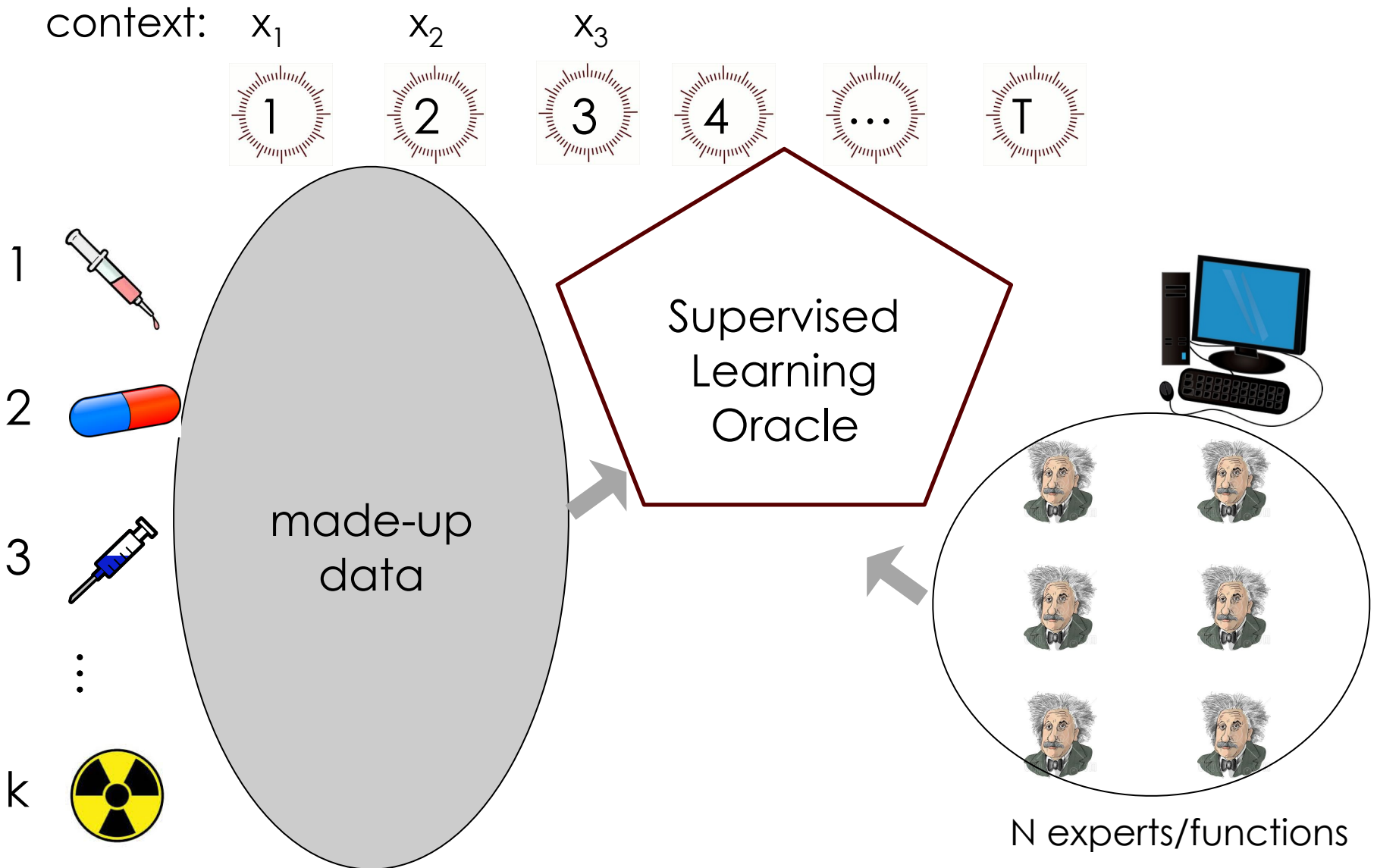


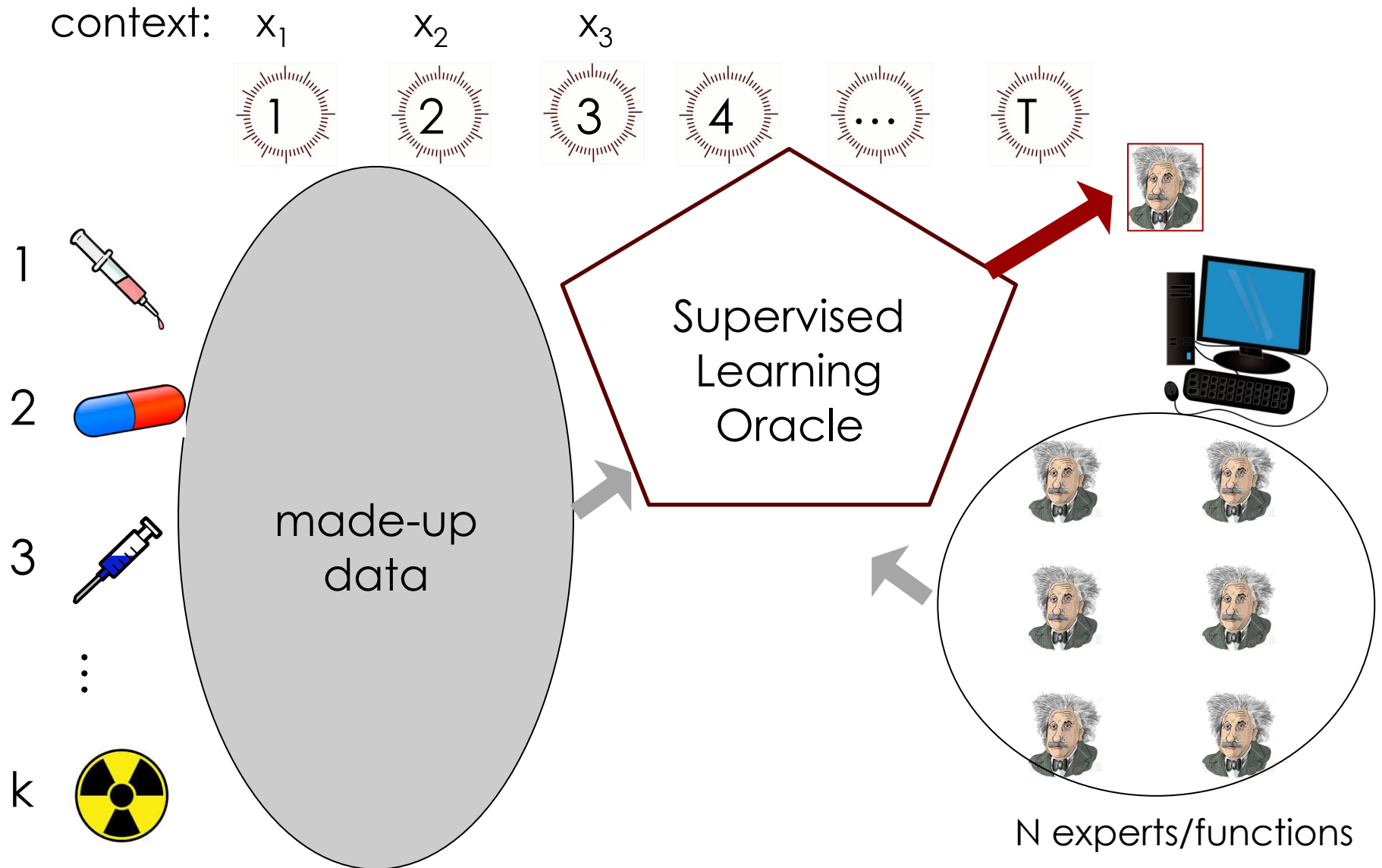
k



3

N experts/functions





**Thm:** [Dudik-Hsu-Kale-Karampatziakis-Langford-R-Zhang '11]:  
For any  $\delta > 0$ , w.p. at least  $1 - \delta$ , given access to a learning oracle, R-UCB has regret  $O((KT \ln(NT/\delta))^{1/2})$  in the stochastic contextual bandit setting and runs in time  $\text{poly}(K, T, \ln N)$ .

Main idea:

make a convex program that optimally  
“solves” the bandit problem.



(Ab) use the supervised learning oracle to act  
as a separation oracle for this problem.

Taming the Monster: A Fast and Simple Algorithm for  
Contextual Bandits  
[Agarwal-et al. '14]





A research goal: make this work in adversarial model.


# Bandit Slate Problems


## [Kale-R-Schapire '11]

Sci/Tech

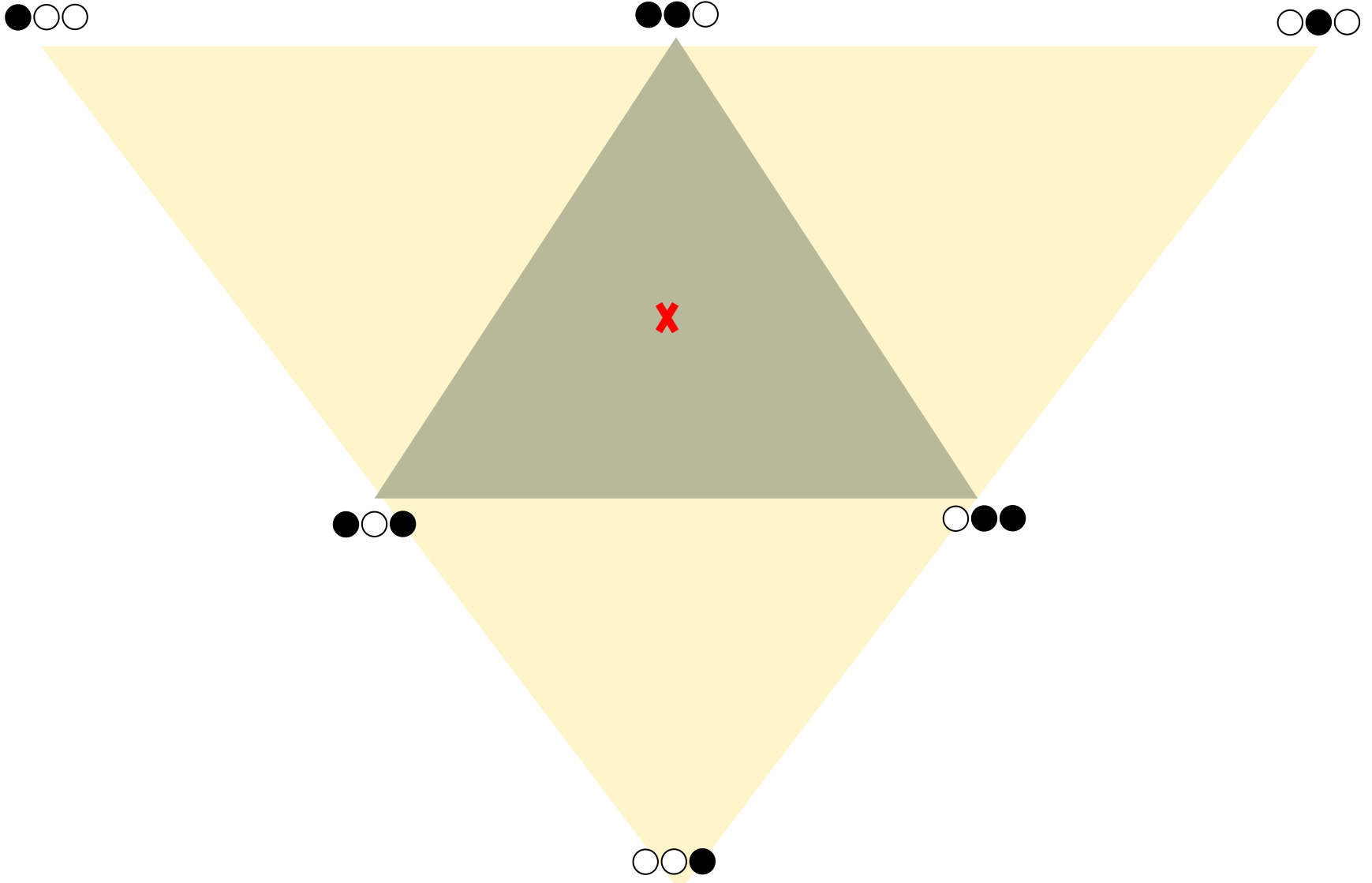
 **LivingSocial hacked; 50 million affected**  
CNET | 1 hour ago | Written by Seth Rosenblatt  
Hackers target LivingSocial, stealing the personal data of more than 50 million people in an enormous security breach. Seth Rosenblatt.

 **Apple loses more global smartphone market-share to Samsung**  
San Jose Merc... | 1 hour ago | Written by Pat May  
Once again, Samsung's smartphone successes have come at Apple's expense. In its latest report, research firm IDC revealed Friday that even though robust iPhone 5 sales helped goose Apple's total smartphone sales in its most recent quarter by 6.6 percent ...

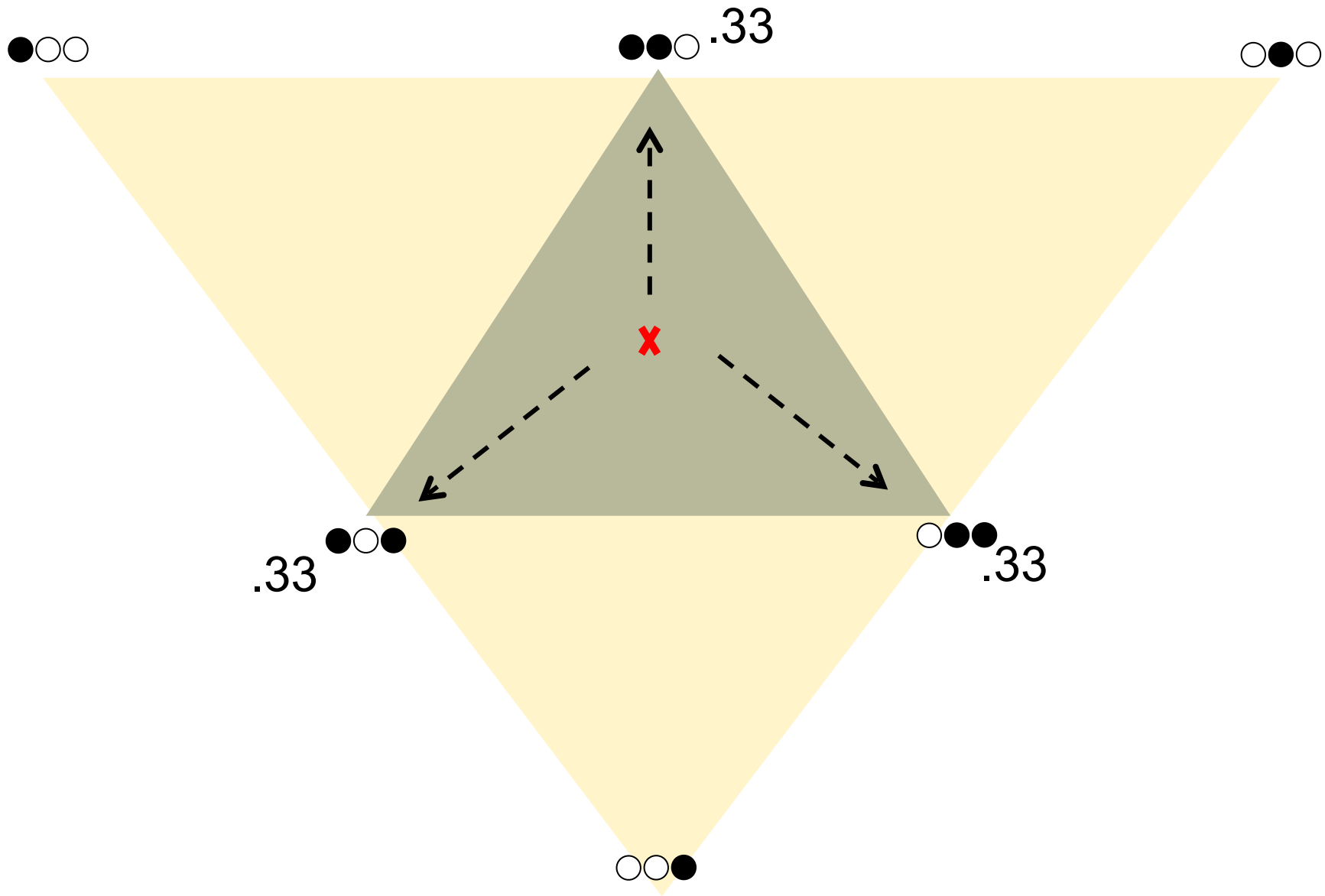
 **Monkeys imitate local food norms, study finds**  
Christian Scien... | 1 hour ago | Written by Mai Ngọc Châu  
The tendency to adapt to cultural behaviors in a new place is not unique to us, a new study suggests. Skip to next paragraph. In Pictures: Monkeying around!

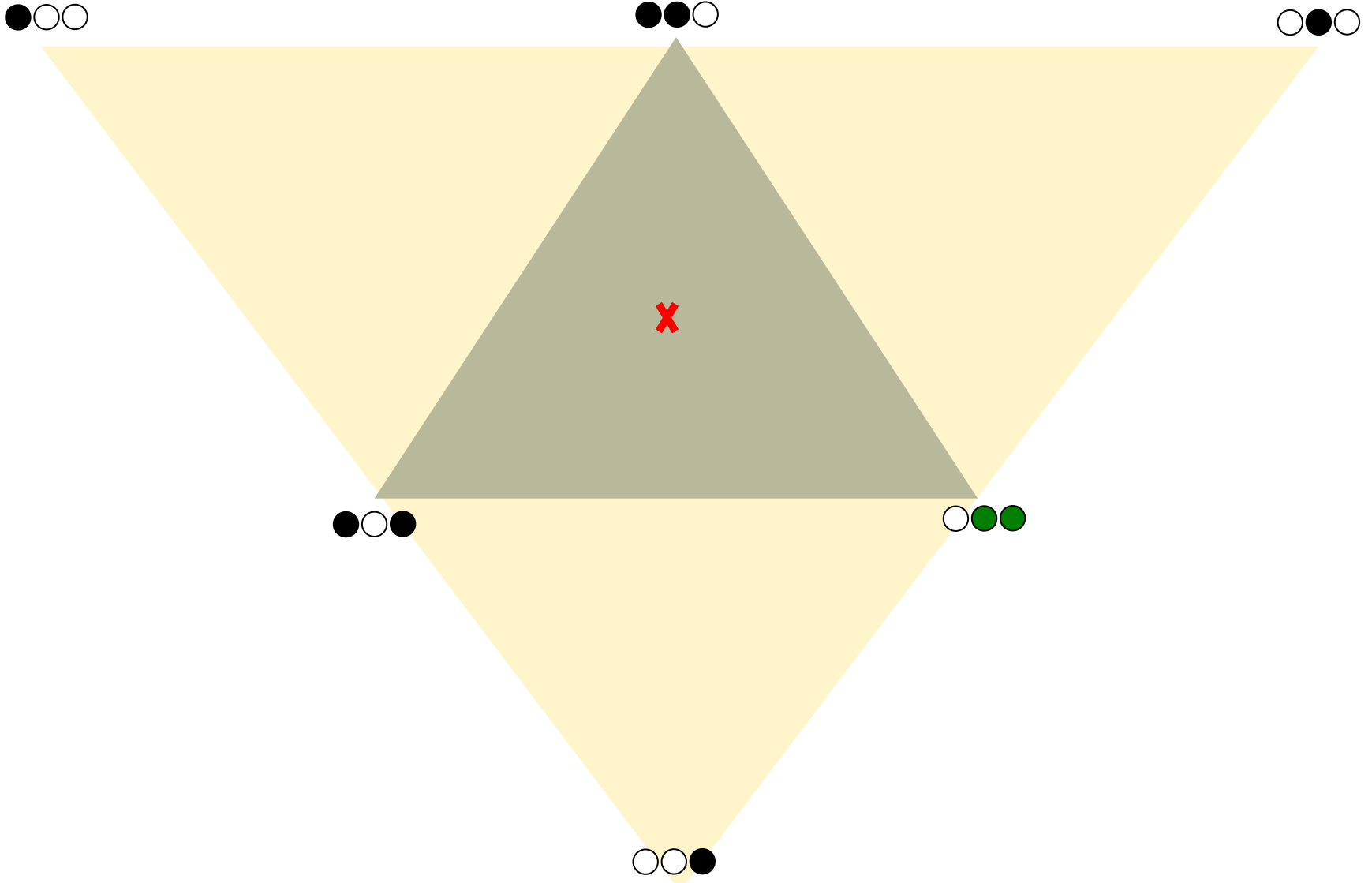
 **Dutch Man Said to Be Held in Powerful Internet Attack**  
New York Times | 1 hour ago | Written by Nicole Perloth  
Dutch authorities say police officials in Spain have arrested a man believed to be connected to an online attack on a spam-fighting site that snarled the Internet last month.

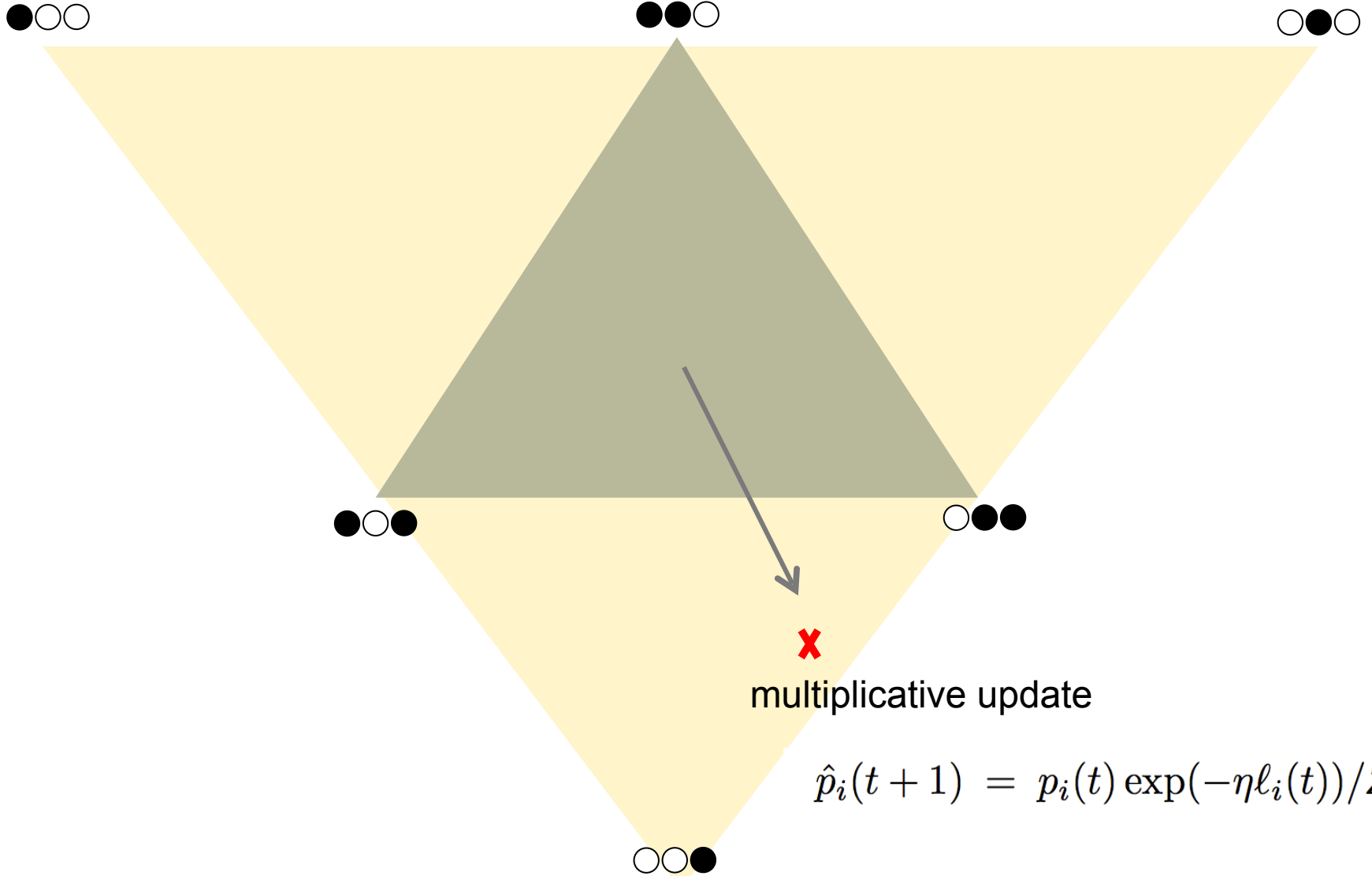
Instead of selecting one article, we need to **select  $s \geq 1$** , articles (possibly ranked). The motivation is web ads where a search engine shows multiple articles at once.





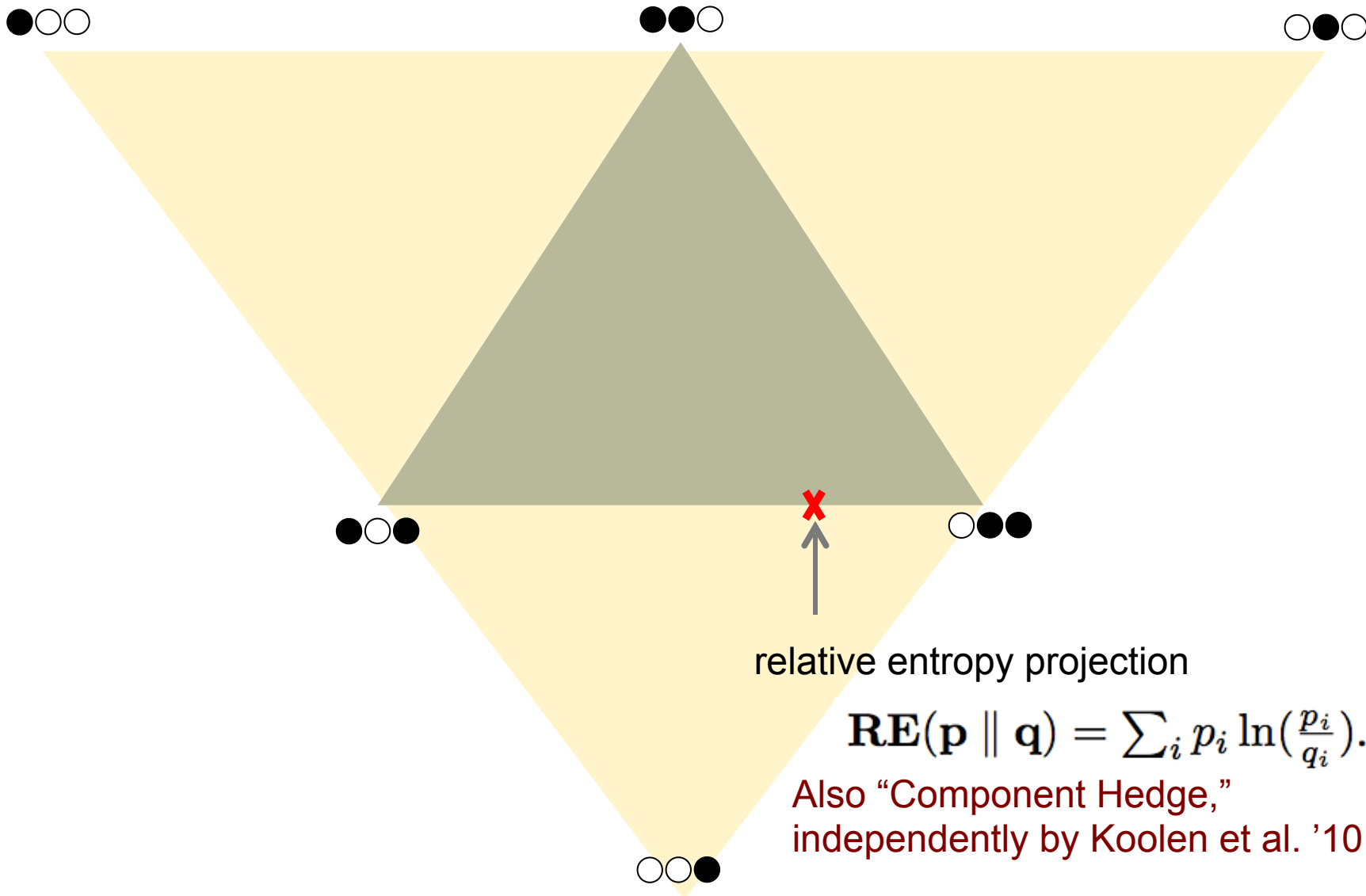


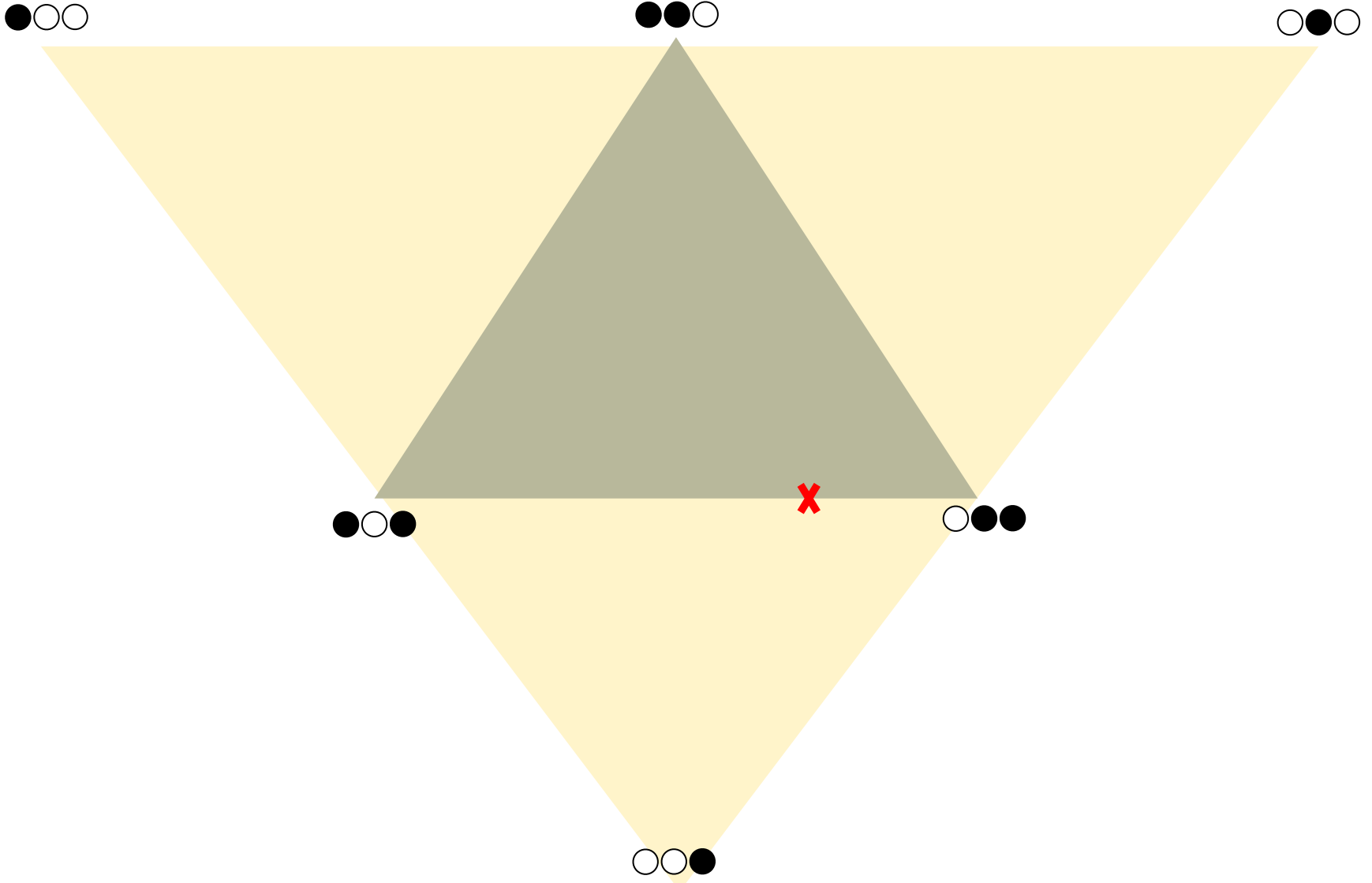




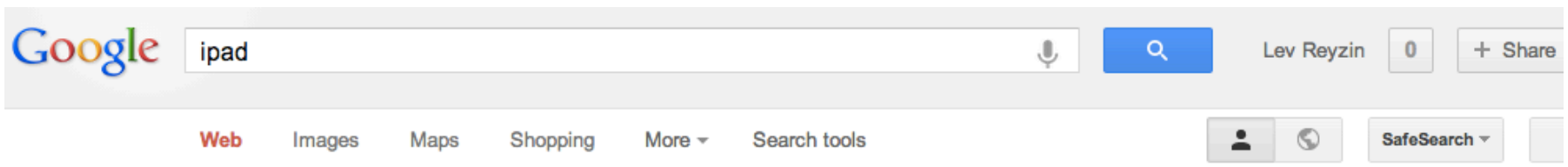
**X**  
multiplicative update

$$\hat{p}_i(t+1) = p_i(t) \exp(-\eta \ell_i(t)) / Z(t)$$





# An Interesting Issue: Ads



160 personal results. 1,810,000,000 other results.

Ad related to **ipad** ⓘ

[Official Apple® Store - Order iPad with Retina display now](http://store.apple.com/ipad)  
[store.apple.com/ipad](http://store.apple.com/ipad)

★★★★★ 209 reviews for store.apple.com  
Free shipping and free engraving.

241 people +1'd this page

10 hour battery life - Free In-Store Pickup - Choose Cellular or Wi-Fi - iPad 2

[Shop for ipad on Google](#)

Sponsored ⓘ



[Apple iPad with Retina...](#)  
**\$479.99**  
AT&T



[iPad with Retina displa...](#)  
**\$499.00**  
Apple Store



[Apple iPad 2 Mc954ll/a Wif...](#)  
**\$389.00**  
eBay



[Apple iPad mini 16GB wit...](#)  
**\$299.00**  
Walmart



[Apple - Ipad Mini...](#)  
**\$329.99**  
Best Buy

Shop by cellular connectivity: [Wi-Fi Only](#) [3G](#)

Ads ⓘ

[iPad at Walmart](#)

[www.walmart.com/Apple\\_iPad](http://www.walmart.com/Apple_iPad)  
Save On iPad at Walmart.  
Free Shipping Site to Store.

[iPad Apple at Amazon](#)

[www.amazon.com/iPad+Apple](http://www.amazon.com/iPad+Apple)  
★★★★★ 736 reviews for amazon.com  
Save Big on New Gear at Amazon!  
Free 2-Day Shipping w/Amazon Prime.

[50% Off iPad](#)

[www.comparedstores.com/iPad](http://www.comparedstores.com/iPad)  
iPad Huge Discounts  
2012 Clearance Sale, Free Shipping!

[Overstock iPads Sale](#)

[www.dealday.com/iPad](http://www.dealday.com/iPad)  
Want Apple iPads At Cheap Prices?  
Signup And Save On All Apple iPads!

[See your ad here »](#)

[Apple - iPad](#)

[www.apple.com/ipad/](http://www.apple.com/ipad/) ▾

iPad is a magical window where nothing comes between you and what you love. And it comes in two sizes.

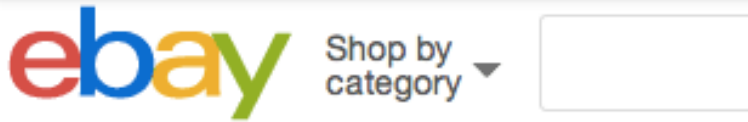
[Shop iPad](#)

iPad 2 - Compare iPads - iPad mini -

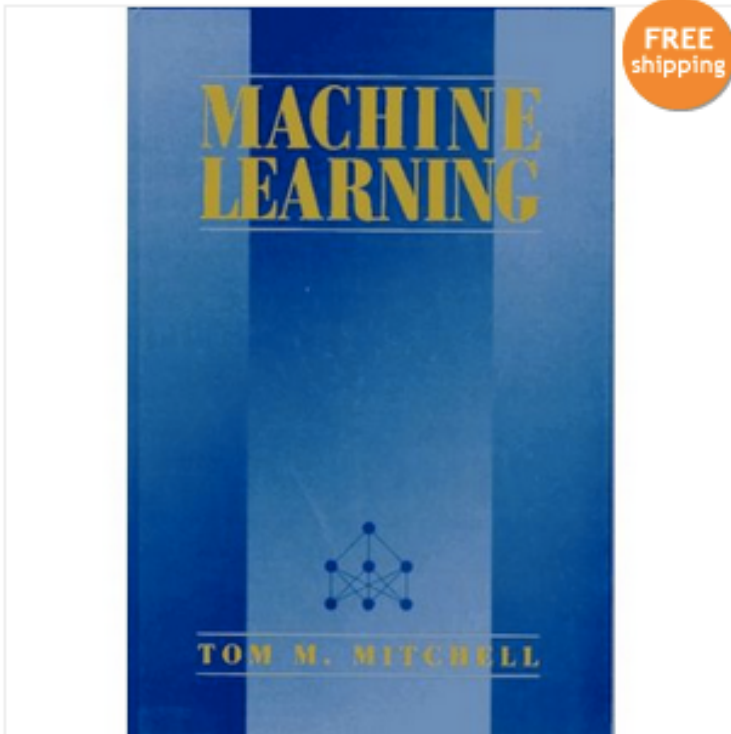
[Compare iPad models.](#)

Whether you choose iPad mini, iPad

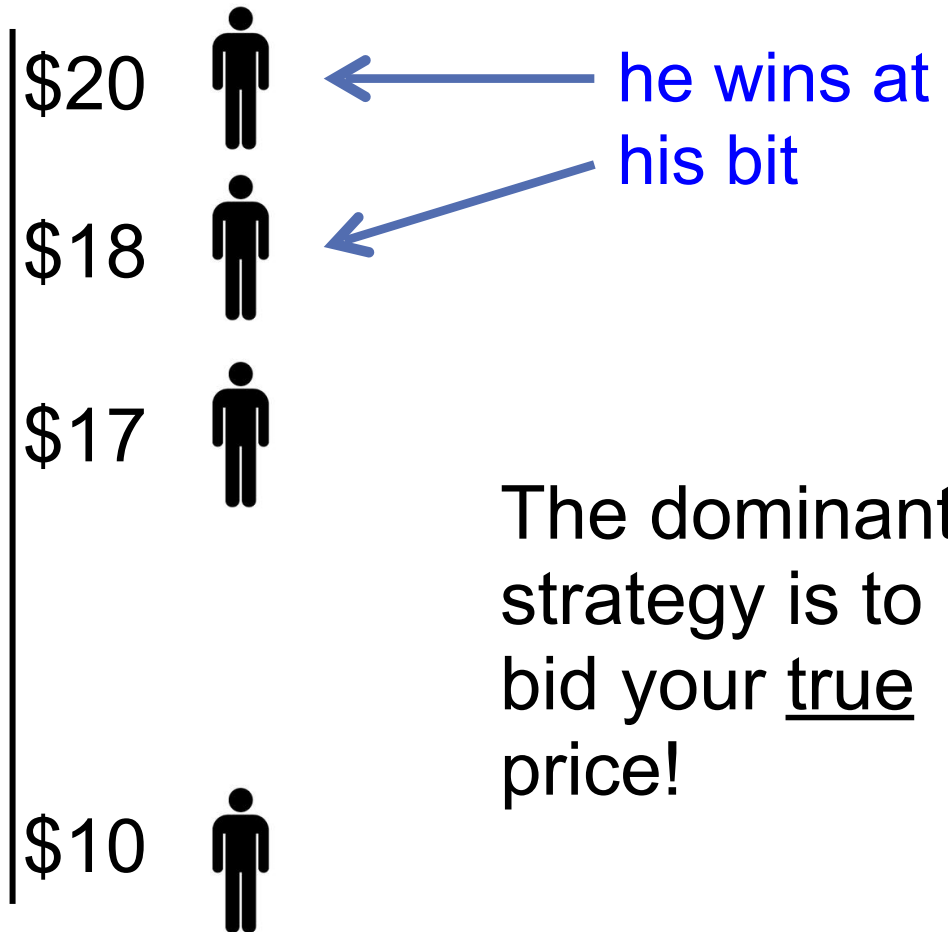
# SECOND PRICE, TRUTHFUL BIDDING



[Back to search results](#) | Listed as Machine Learning t



Click to view larger image



The dominant strategy is to bid your true price!

# CONTEXTUAL ADVERTISING

1. **For an ad to be shown, it must have high **expected earnings**.**  
Earnings = clickthrough rate (CTR) x expected charged price
2. **CTR must be **learned****  
a classic contextual bandits problem
3. **Charged prices are functions of the bids of advertisers.**  
e.g. Can't ever charge more than an advertiser's bid
4. **Ads must be shown so that CTR is learned quickly, but the auction should be **truthful**.**



# **SUMMARY**

**When dealing with many customers/ subscribers and many options, a smart automated strategy needs to be employed.**

**This is becoming true of nearly every company presenting content online.**

**Presents many important mathematical challenges, most of which are wide open.**

**THANK YOU!**