# Motivating problem

- Each agent receives a piece of data in rounds $t = 1, \ldots, T$
  *e.g. a bit*
- Accuracy goal: maintain accurate statistics
  *e.g. average of agents' bits at the current time*
- Privacy guarantee: over **entire** time horizon in **local** model
  *agents hold their data; all communications they make are d.p.*

**One approach:** randomized response on each $t$ separately.
*Accuracy degrades polynomially in $T$*

If data changes are **slow** or **rare**, we hope to do better!

# Our setting

**Stochastic setting:** data is drawn from a distribution.
**Assumption 1:** users all draw from the same distribution!
*Agent $i$'s bit at time $t$ is $b_i^t \sim Bernoulli(p^t)$*

**Examples:** auditing gambling systems, product defect rates.
*(contrived?)*

**Assumption 2:** distributions change only $k$ times out of $T$ rounds.
$\implies$ *for fixed $\varepsilon$, accuracy "should" only degrade with $k$.*

# Our approach

**(1)** split rounds into **epochs**:

- Within an epoch, users aggregate their own data.
  → *obtains estimate of distribution during that epoch*
- After each epoch, users report to the center.
- Center publishes accurate statistics after each epoch.

**(2)** Use a **consensus** protocol to detect changes:

- Users who detect a significant change in distribution vote YES
  *using randomized response*
- If enough YES votes, center initiates a **global update**
  *estimated distributions are reported and aggregated with RR*
- W.h.prob, agents only vote/update $\Theta(k)$ times

# Key technical challenge

If a **small** change occurs:

- Accuracy is not affected...     *if it were, an update would trigger*
- ...but privacy may be!     *YES voters are repeatedly ignored*

**Solution: synchronized intensity-frequency protocol**.
If you detect a _____ significant change, vote YES, but only if...

- **very**: always vote YES.
- **somewhat less**: only if $t \mod 2 = 0$.
- **even less**: only if $t \mod 4 = 0$.
- **...**: only if $t \mod 2^\ell = 0$.
- **...almost insignificant:** only if $t = 0$ or $T/2$.

# Why does it work?

If you detect a **very** significant change, you can be confident. . .

- **not:** many others also did. . .
- **but:** many others detected a **somewhat less** significant change!
- $\implies$ by the time you vote twice, a vote will succeed.
- Once a vote succeeds, a global update occurs
  $k$ *changes* $\implies O(k)$ *YES votes and updates*

*Less-frequent turtles all the way down!*

# Results summary

## Theorem (Privacy)

*Each user is guaranteed $\varepsilon$-local differential privacy.*

*Holds without any assumptions.*

## Theorem (Accuracy)

*With high probability, when epochs are of length $\ell$ and $n$ users, global estimate of $p^t$ is accurate to $\frac{k \log T}{\epsilon} \left( \frac{1}{\sqrt{\ell}} + \frac{1}{\sqrt{n}} \right)$.*

*Under assumptions on same distribution and $k$ changes.*

# Extensions and Directions

- Extension: histograms
  *Can integrate with e.g. Bassily-Smith 2015; more work needed*
- Extension: multiple subpopulations
  *as long as each has $\geq \sqrt{n}$ members*
- Direction: other algorithmic approaches
- Direction: other models
- Direction: lower bounds

**Thanks!**