

Survivability of Cloud Databases Factors and Prediction

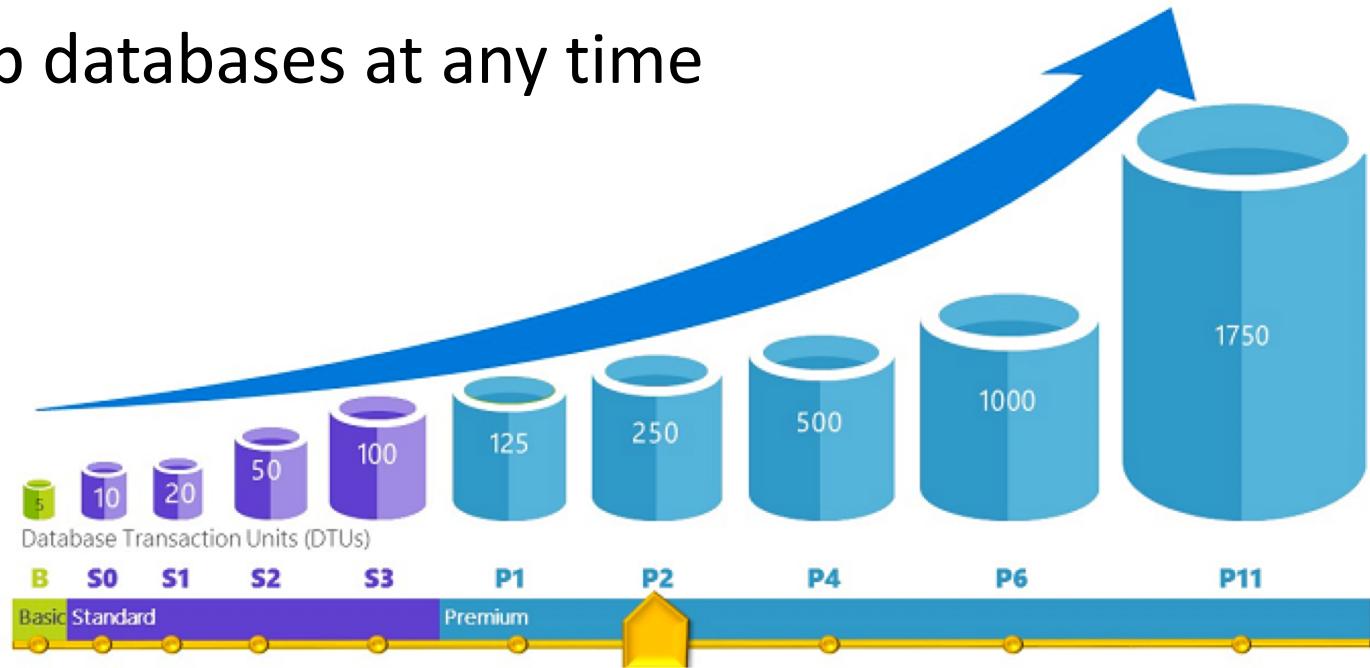
Jose Picado

Willis Lang

Edward C. Thayer

Microsoft Azure SQL Database

- Microsoft's cloud relational database service: Azure SQL Database
- Azure SQL Database offers different service editions
 - Basic, Standard, Premium
- Users may create or drop databases at any time



Life cycle of a database

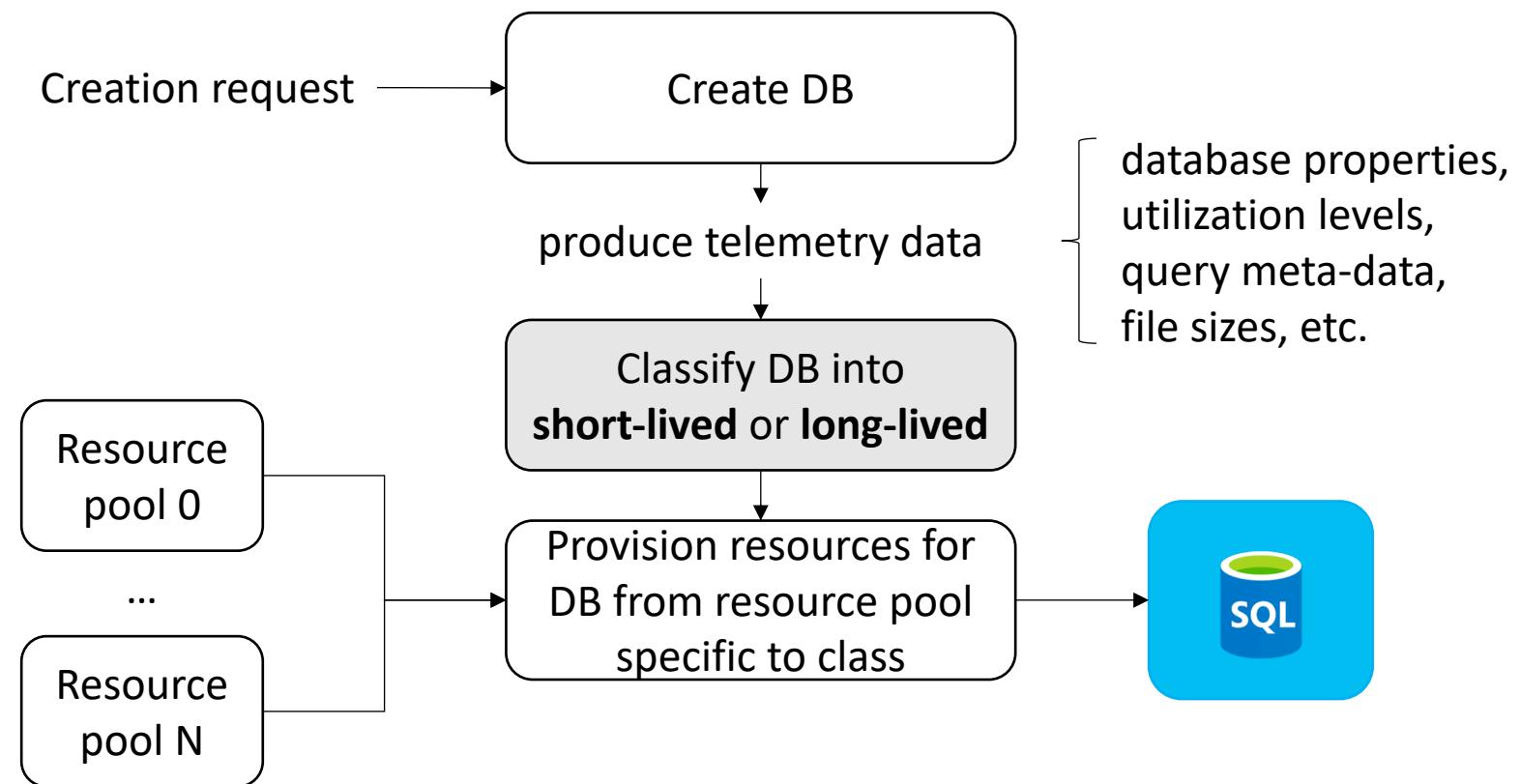
- The life cycle of a database consists of:



- During this life cycle:
 - Use can change the database's edition or scale up and down within edition
 - Database produces telemetry data
 - Resources need to be provisioned for the database

Longevity-based computer resource provisioning

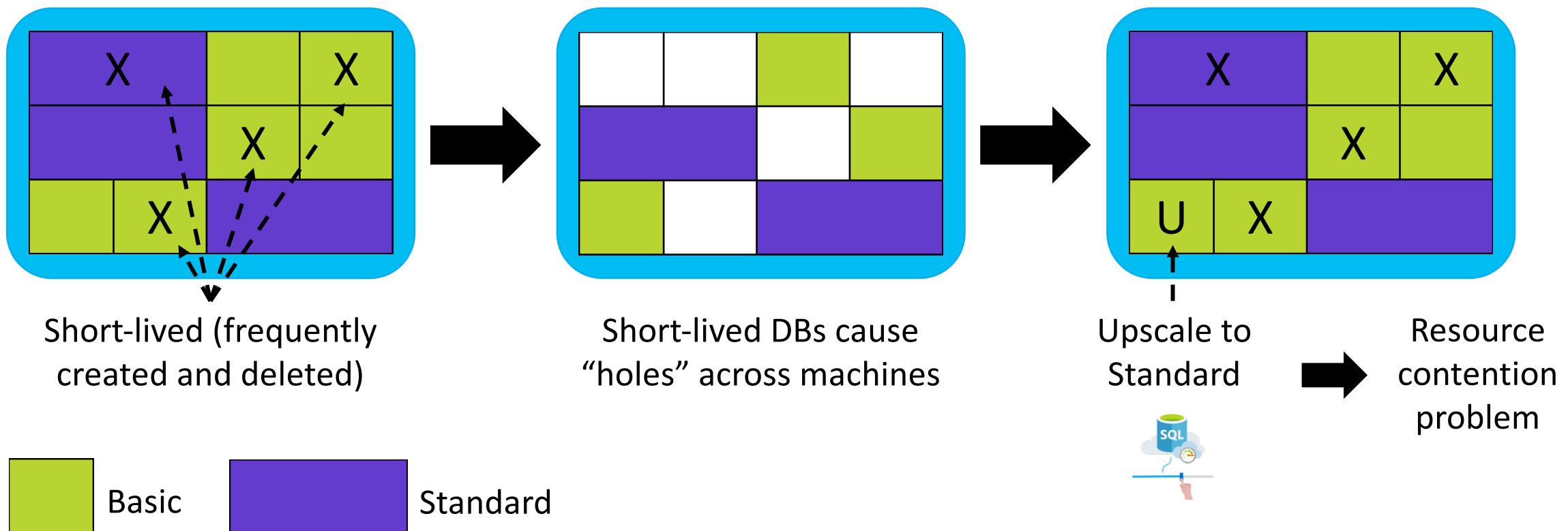
- Provision resources to databases based on its predicted longevity



Why smart provision?

Alleviate resource contention problems

- Frequent creation and deletion of short-lived databases cause resource contention problems

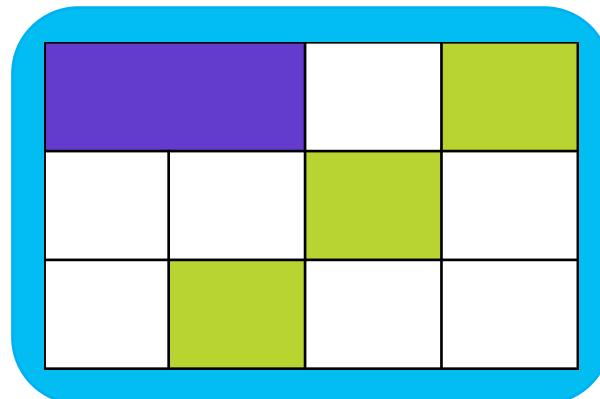


Why smart provision?

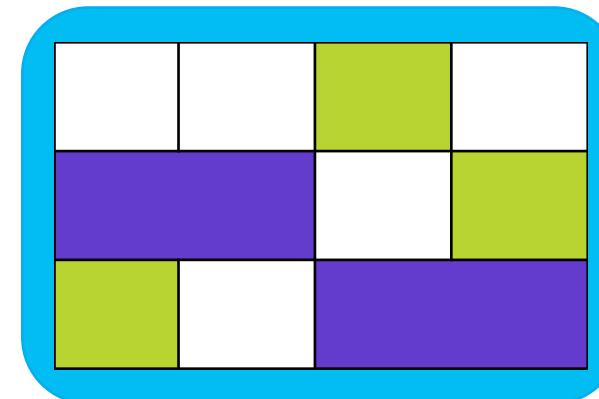
Alleviate resource contention problems

Solutions:

1. Set aside extra capacity in case databases are upscaled – wasteful!
2. Defragmentation – expensive!
3. Smart separation of short-lived and long-lived databases



Short-lived databases

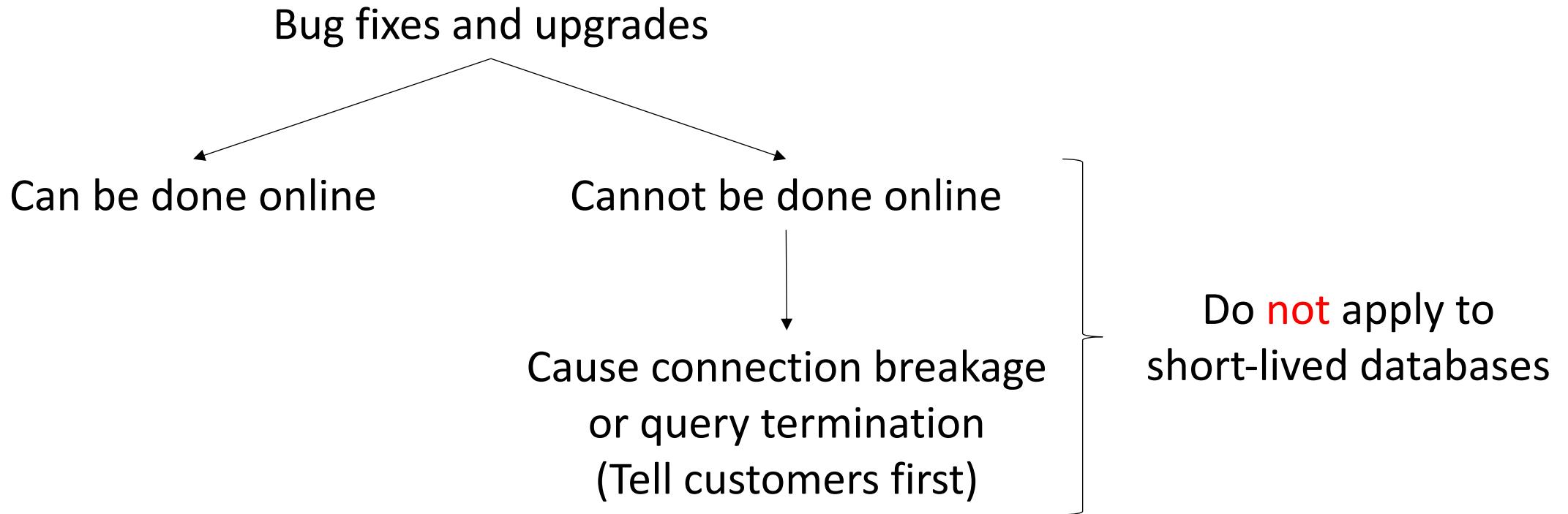


Long-lived databases

Why smart provision?

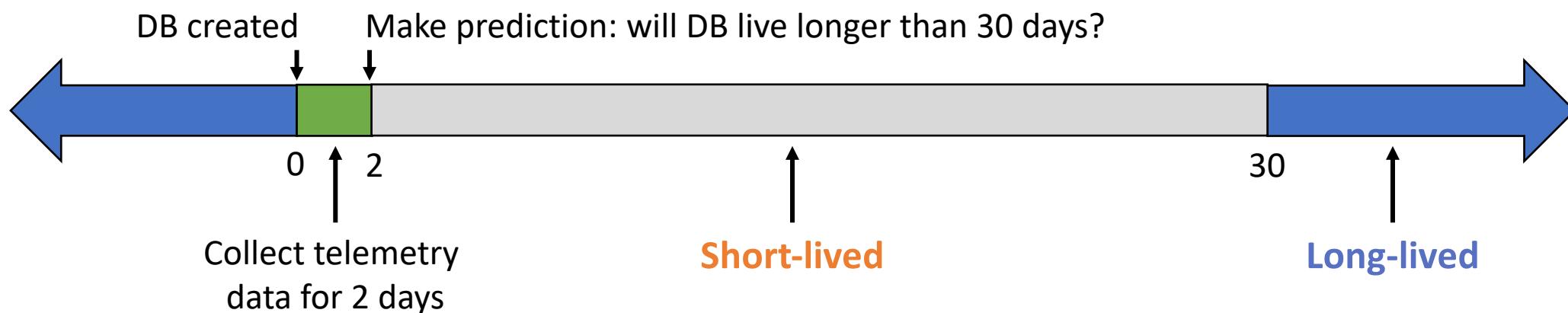
Avoid service disruptions

- Short-lived databases have short life cycle



Classification problem

- Classify into **short-lived** ($\text{live} \leq 30$ days) or **long-lived** ($\text{live} > 30$ days)
- Collect telemetry data for first 2 days after creation
- Extract features: creation time, database name, database size, subscription type, subscription history, etc.



Experimental setup

- Dataset: Azure SQL databases created in a period of five months
- Three full production regions around the world
- Basic, Standard and Premium editions
- Databases that live ≥ 2 days
 - Some subscriptions create *only* databases that live less than 2 days
 - Usage pattern: frequent cycling of databases
- Classifier: random forests



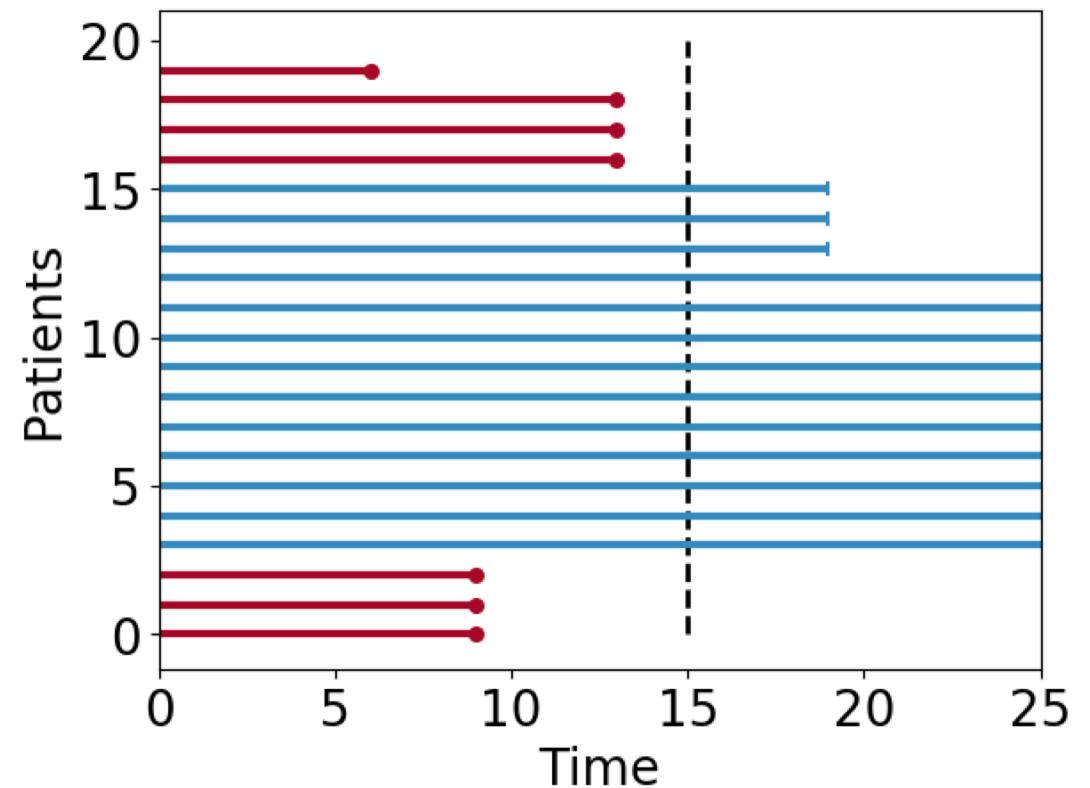
Evaluate predictions using survival analysis

- **Survival analysis:** analyze the expected duration of time until an event happens
- Used by actuaries and medical researchers to measure lifetimes of populations
- We leverage well-accepted statistics to evaluate our model and to compare survival distributions of short-lived and long-lived databases

Concept	Our case
patient	databases
birth	database is created
death	database is dropped
treatment	assign class using our model and provision resources specific to class

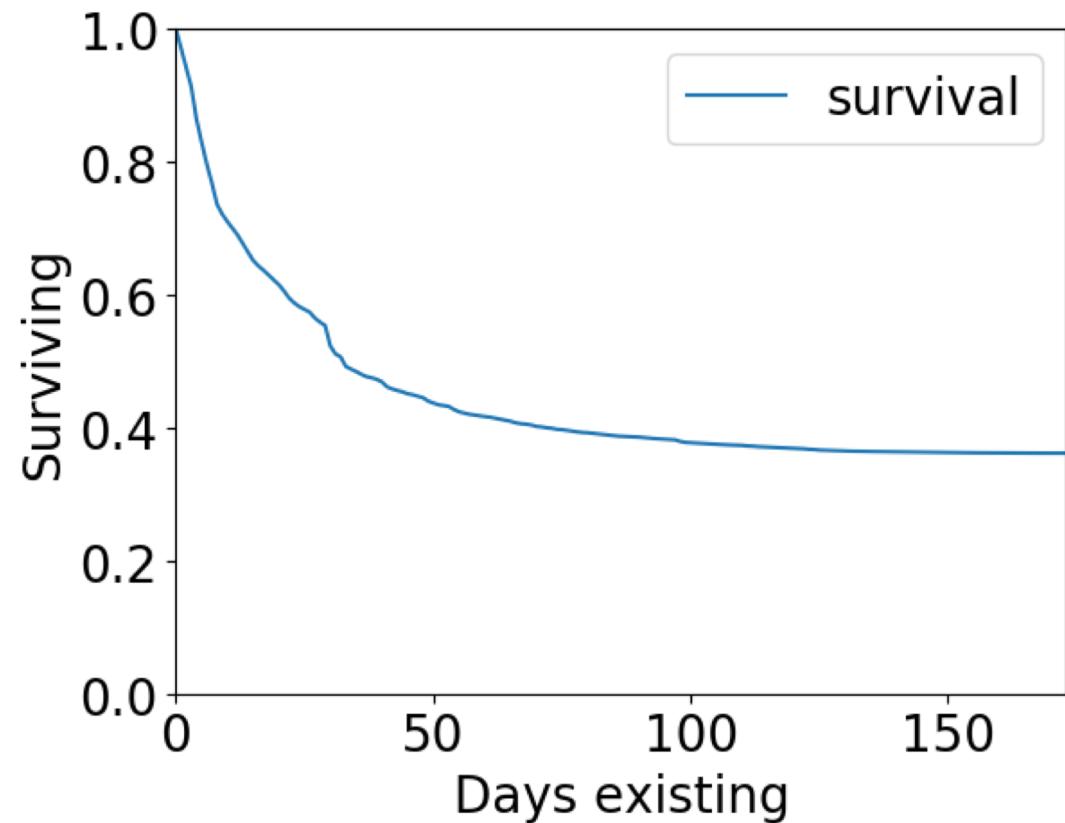
Censored data

- We haven't seen the death event for some individuals in the population yet (some databases have not been dropped yet)
- These individuals are **censored**
- Want to use this data without being biased by censorship



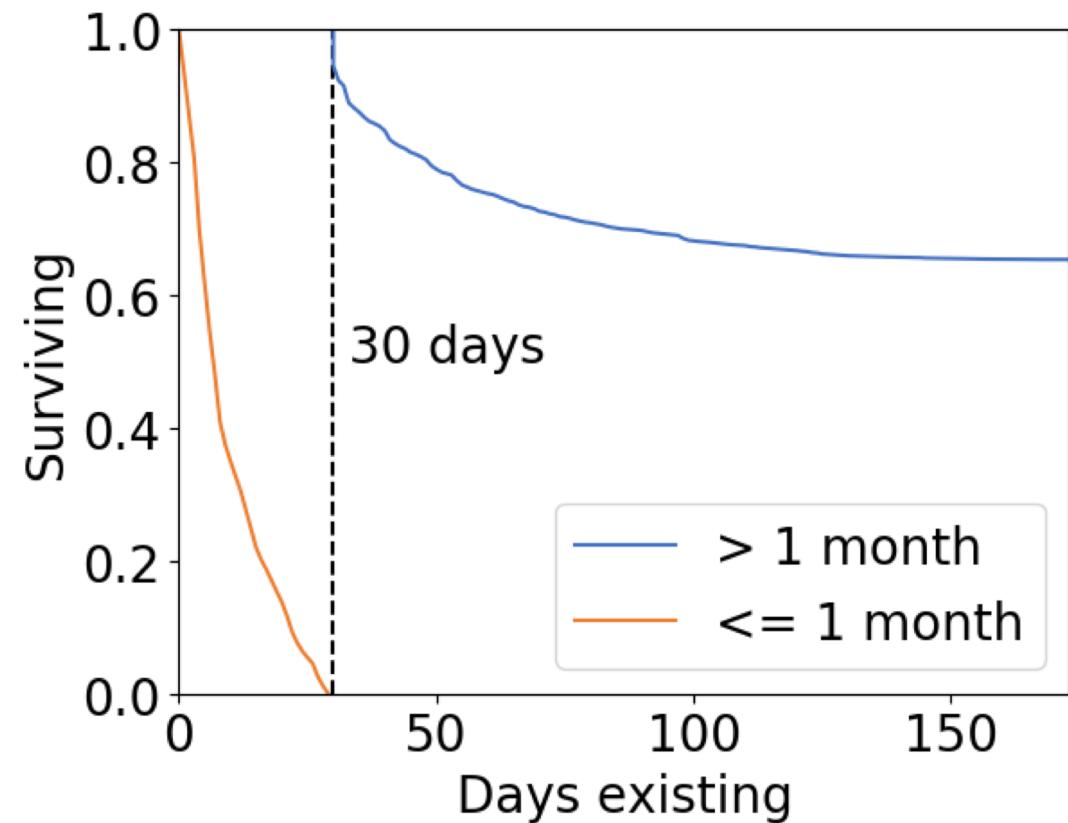
Survival curve of Azure SQL DBs

- **Kaplan-Meier survival curve** shows the probability that an individual lives longer than time t
- E.g., there's a 40% probability that a database lives longer than 100 days



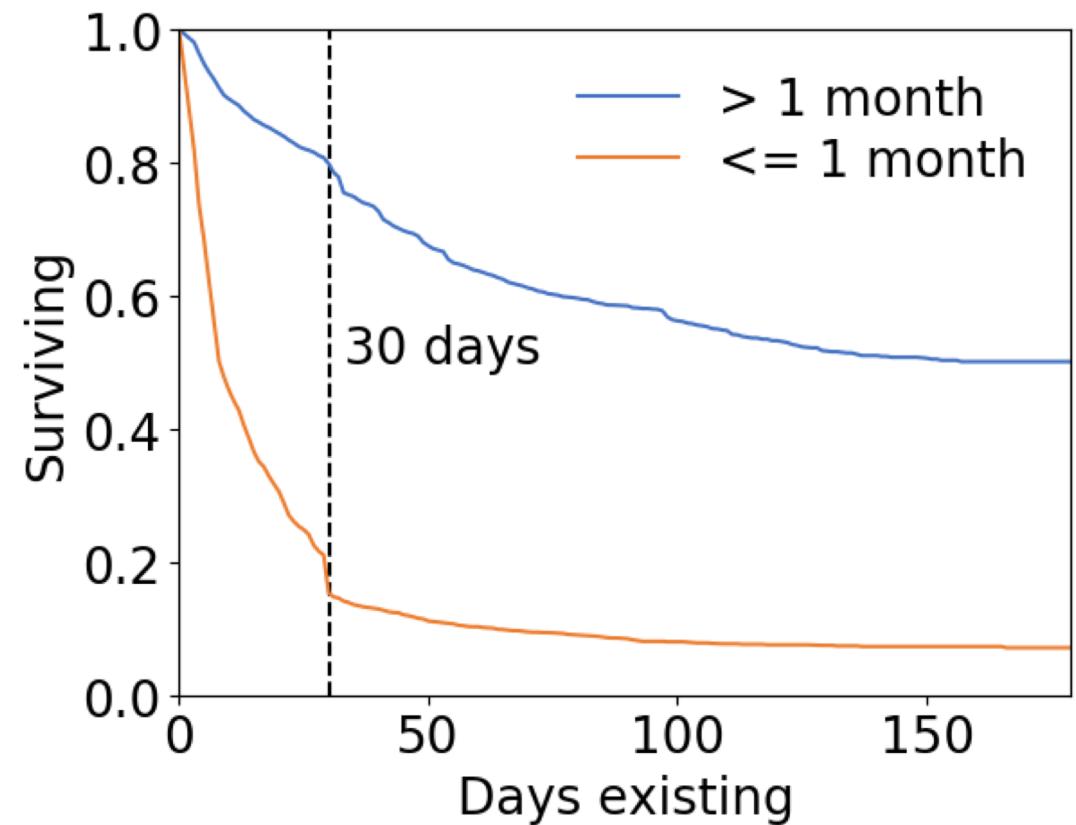
Comparing survival curves

- Separate test DBs by their treatment (class assigned by our model)
- Evaluate predictions by comparing the survival curves of the two groups
- Survival curves for **ideal predictions**:
 - Short-lived should drop to 0% at day 30
 - Long-lived should not fall below 100% until day 31



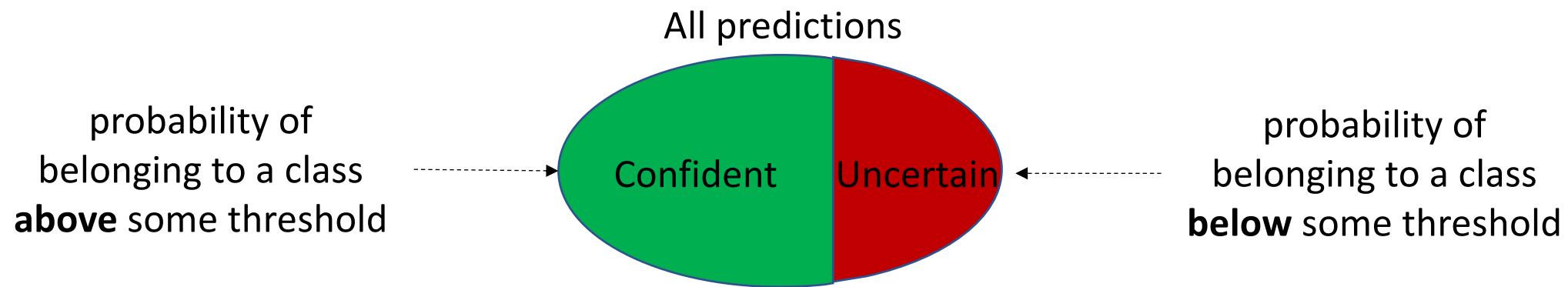
Survival curves for all predictions

- Actual survival curves for all predictions on test databases



Confident/uncertain predictions

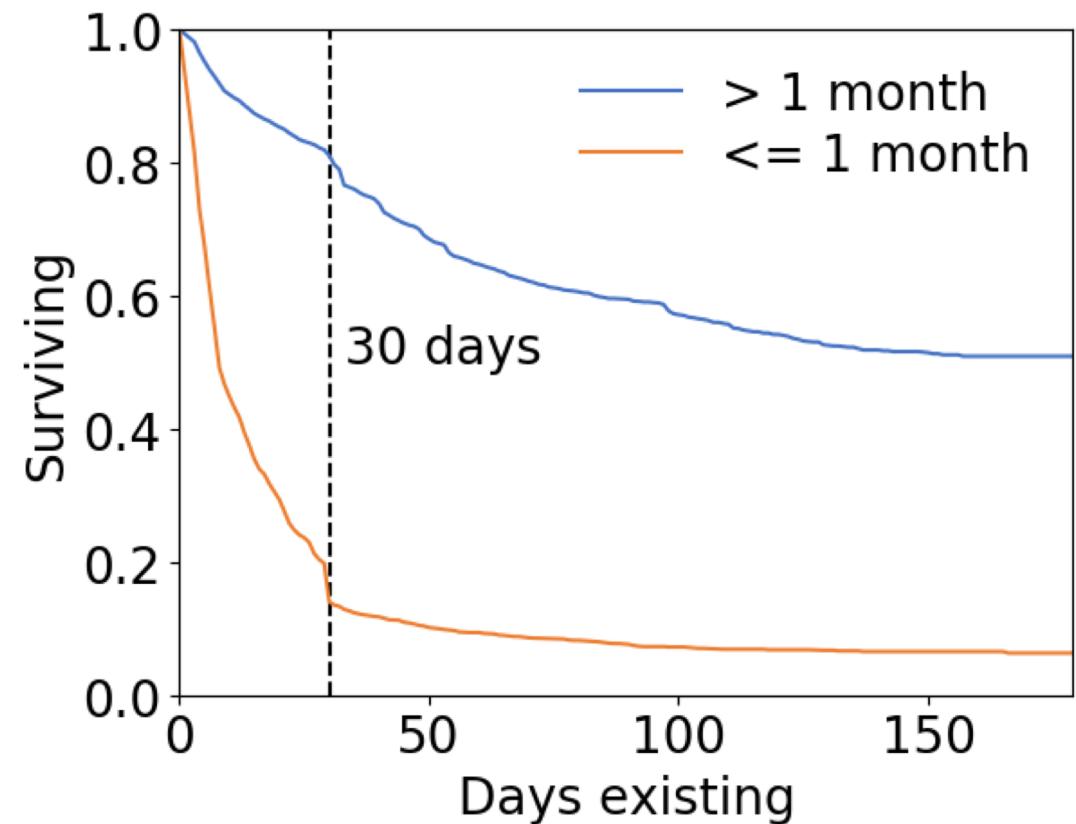
- There may be costs of misclassifications (incorrect resource provisioning decisions)
- Our model can determine how confident it is on its predictions
- Two types of predictions on the test set: confident and uncertain



Survival curves for confident predictions

- **Confident** predictions only
- **Log-rank test**: compare the survival distributions of two groups
- Separation of classes is statistically significant by log-rank test

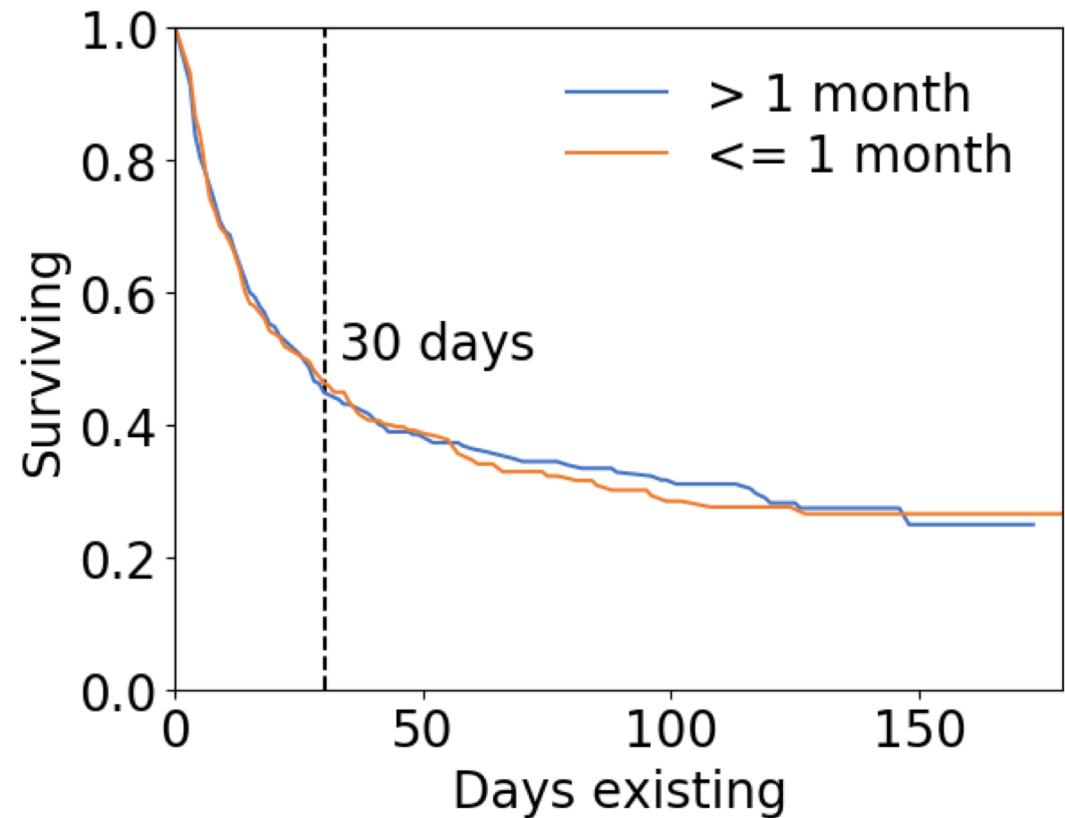
Edition	Confident	Uncertain
Basic	63%	37%
Standard	90%	10%
Premium	71%	29%



Survival curves for uncertain predictions

- **Uncertain** predictions only
- Separation of classes is **not** statistically significant by log-rank test
- If my prediction for a database is uncertain, do not apply treatment to this database

Edition	Confident	Uncertain
Basic	63%	37%
Standard	90%	10%
Premium	71%	29%



Most predictive features

Importance	Category	Features
1	Subscription history	number of previous databases statistics about lifespan
2	Server and database name	length distinct character rate
3	Creation time	hour of the day day of the week week of the year

Conclusions and future work

- It is not trivial to predict lifespan of databases
 - Our model is able to separate short-lived and long-lived databases
 - The difference of survival distributions of databases in each class is statistically significant
- Model is able to determine in which predictions it is confident
 - Reduce misclassification costs by only applying treatment to confident predictions
- Results can be improved with more data (every day new data is generated)
- Explore other features: utilization, logins, sessions, queries

Thanks

Backup

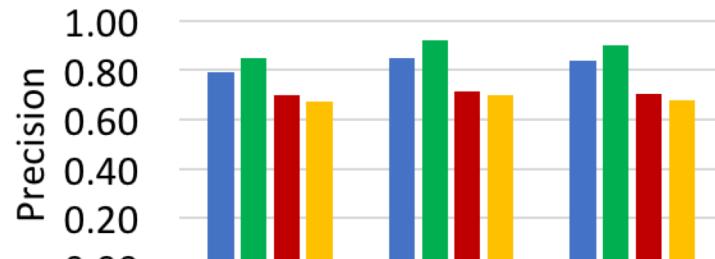
Features

- Create a set of features from raw telemetry data
- Detailed explanation of features can be found in the paper

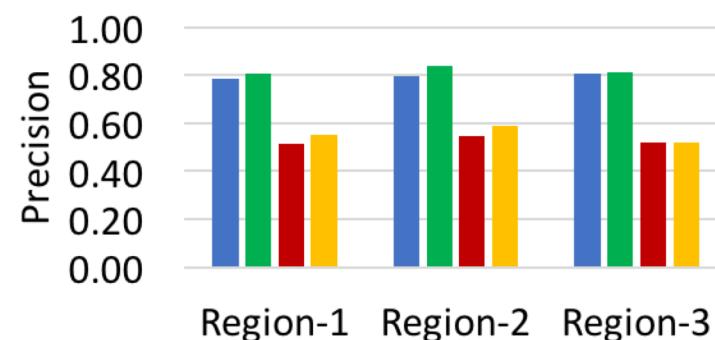
Category	Features
Creation time	hour of the day, day of the week, week of the year, month of the year, ...
Server and database name	length, distinct characters, letter and digits, uppercase and lowercase, ...
Database size	statistics about absolute database size, change in size, ...
Edition and performance level	number of changes, edition/performance level at time of prediction, ...
Subscription type	one-hot encoding of subscription types (trial, consumption, benefits, ...)
Subscription history	number of previous databases, statistics about size and lifespan, ...

■ All ■ Confident ■ Uncertain ■ Baseline

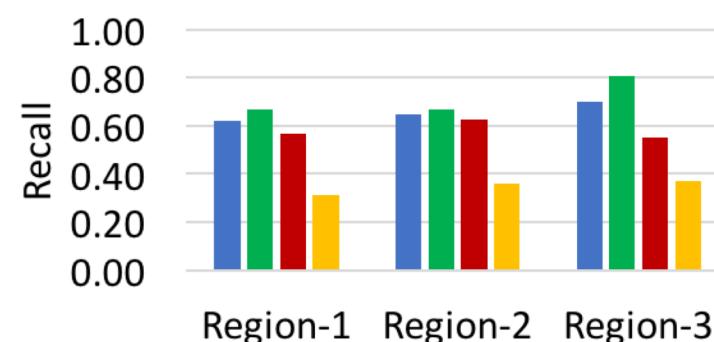
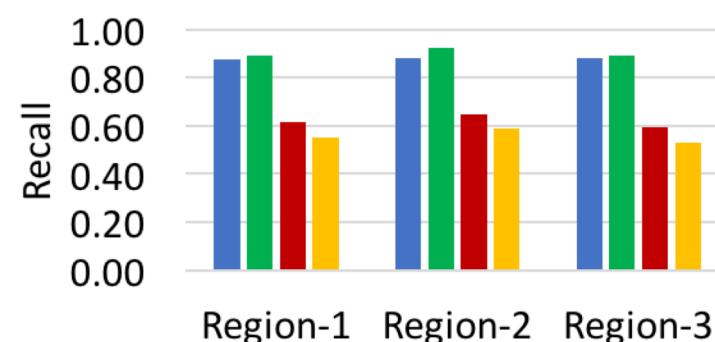
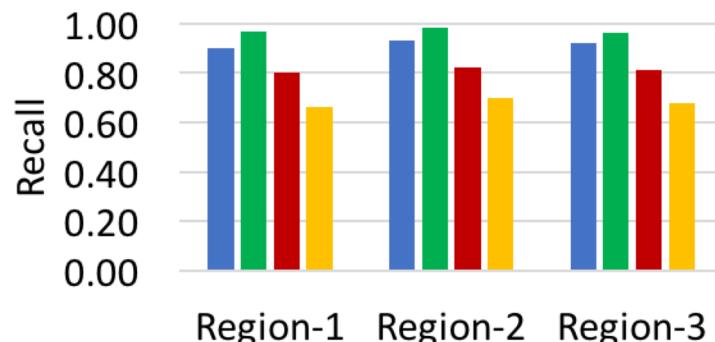
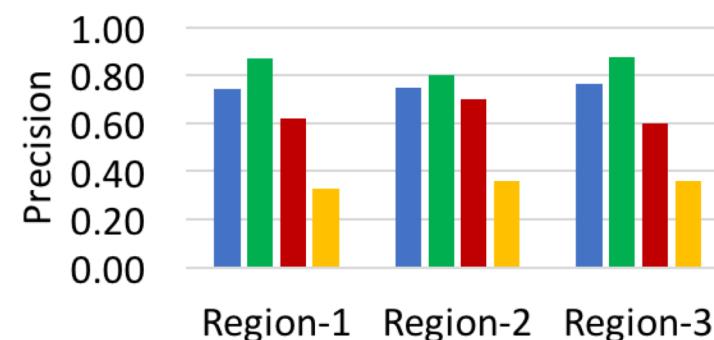
Basic



Standard



Premium



Region	Confident	Uncertain
Region-1	58%	42%
Region-2	63%	37%
Region-3	68%	32%

Region	Confident	Uncertain
Region-1	92%	8%
Region-2	82%	18%
Region-3	97%	3%

Region	Confident	Uncertain
Region-1	71%	29%
Region-2	69%	31%
Region-3	73%	27%