

Outgrowing Your Laptop With Positron

Austin Dickey / 2025-09-17

posit
conf(2025)





Austin Dickey

austin3dickey

Follow

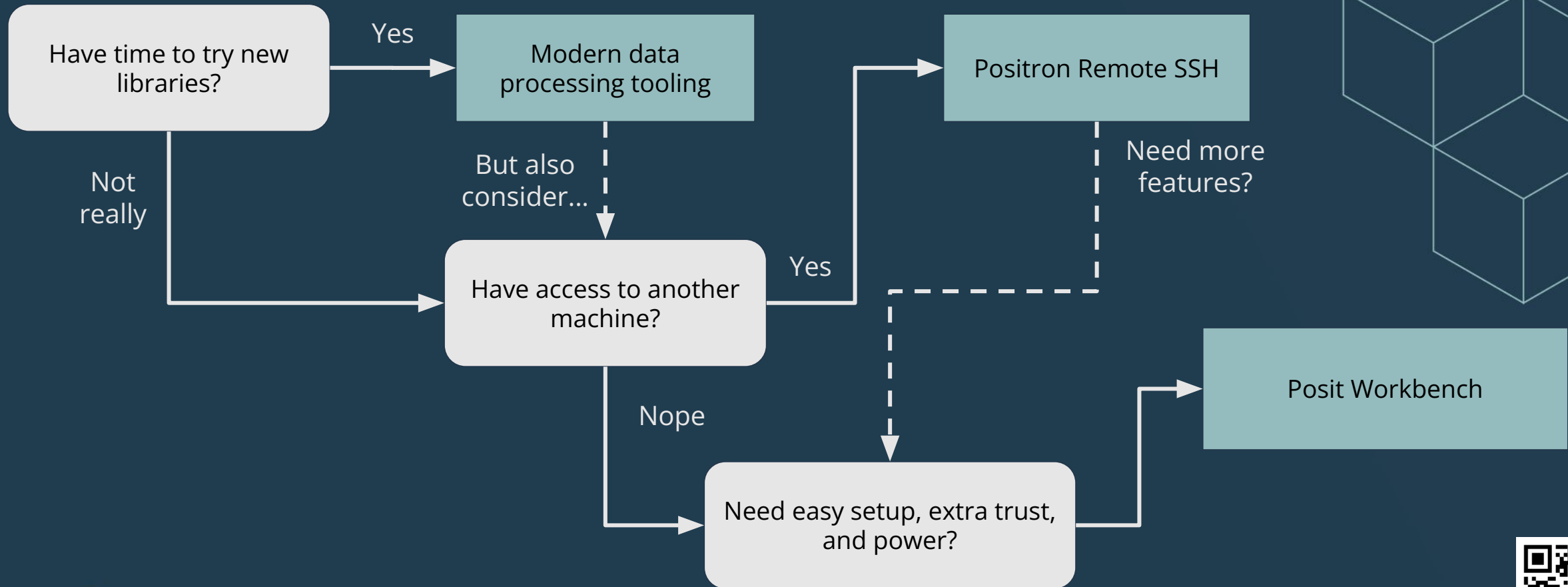
Workin' on Positron at Posit.

posit
conf (2025)

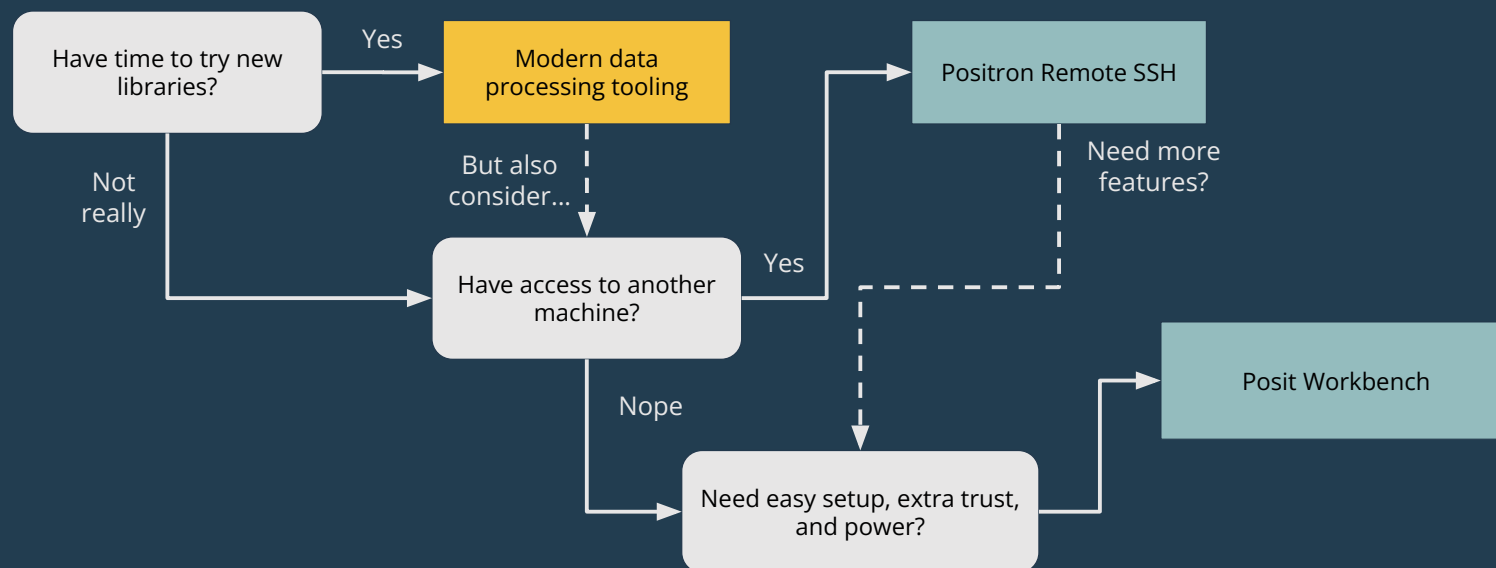
GitHub repo with demo and
slides here



What do I do when I'm struggling to analyze data in memory on my laptop?



Modern Data Processing

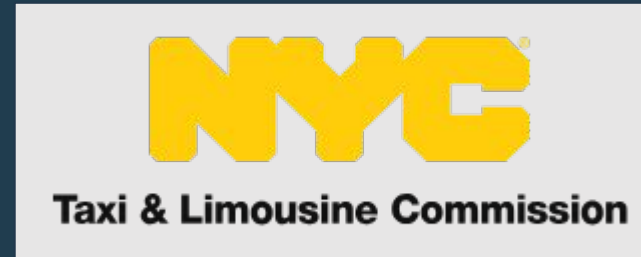


Who is this for?

- People who don't have access to any other machines
- Quick, local iterations
- Proofs of concept
- (Though these techniques are always useful!)



The NYC Taxi Dataset

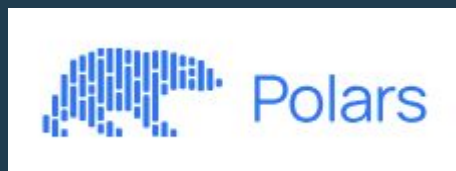


- All trips between 2009-2019
- ~36GB on disk in parquet format
- Would be ~400GB if loaded into pandas



What can I do?

- Luckily there are so many tools that help with this!
 - R and Python work with most of these



What can I do?

- They typically use the following strategies:
 - **Deferred/lazy evaluation:** don't query the data until the last second
 - **Predicate pushdown:** only read the data you need
 - **Split the work** among multiple processors if possible



In Python, I like Ibis, which uses any of these tools.

```
import ibis

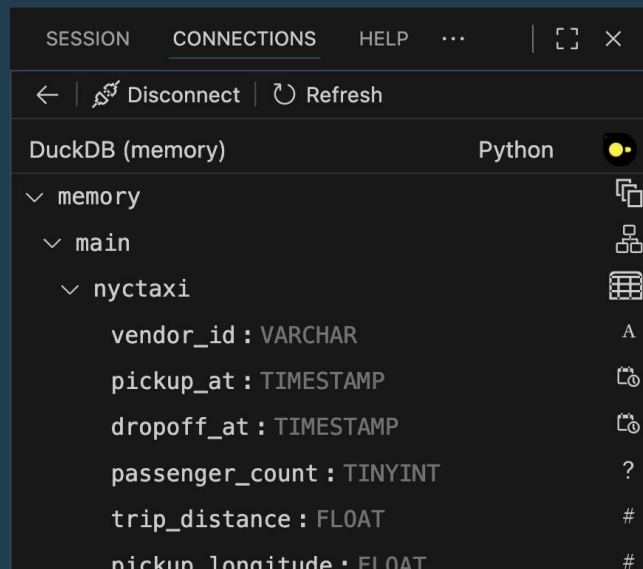
# Set up any backend
con = ibis.duckdb.connect()
con = ibis.polars.connect()
con = ibis.pyspark.connect(session=SparkSession.builder.getOrCreate())
con = ibis.connect(f"postgres://{user}:{password}@{host}:{port}/{database}")

# Load an ibis table (this doesn't actually read the data yet)
t = con.read_parquet(
    "s3://ursa-labs-taxi-data/**",
    table_name="nyctaxi",
    union_by_name=True,
)
```



Why in Positron?

- **Help Pane** lets you learn new tools faster
- **Connections Pane** lets you understand your databases better
- **Data Explorer** lets you understand your data better

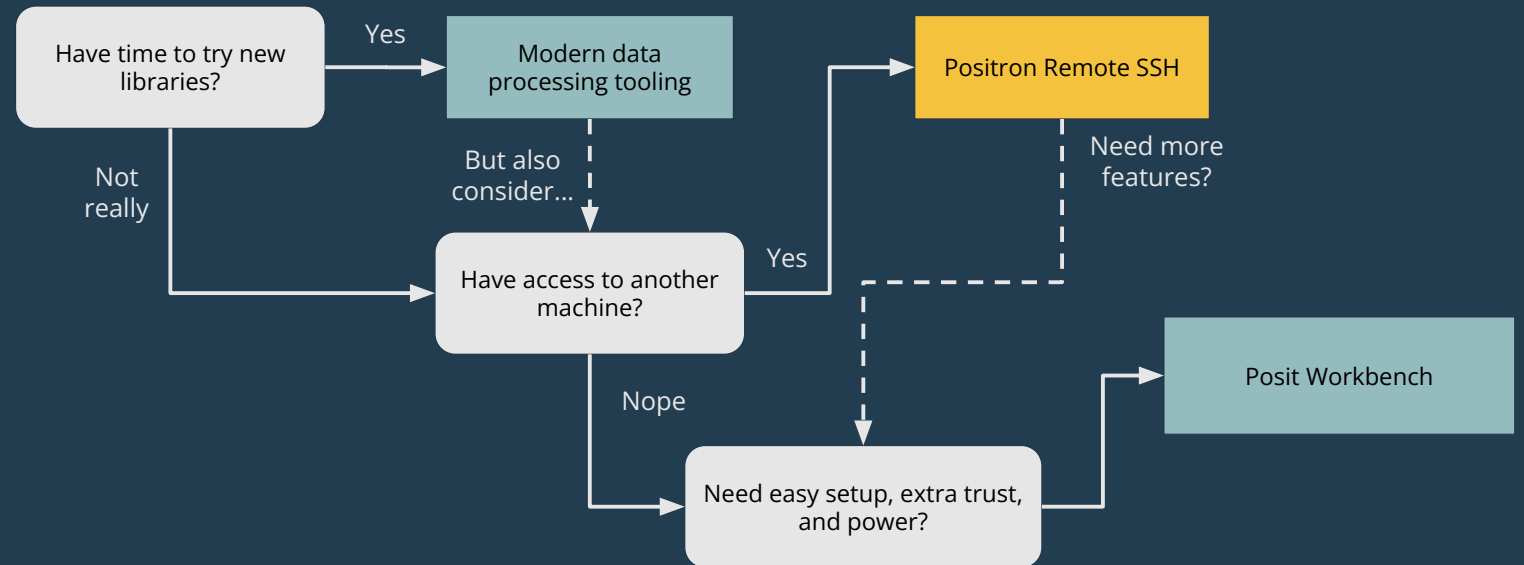


The screenshot shows the 'Data Explorer' pane in Positron, displaying a table of taxi trip data. The table has columns: 'vendor_id' (string), 'pickup_at' (datetime64[us]), and 'dropoff_at' (datetime64[us]). The data is sorted by 'vendor_id'.

	vendor_id	pickup_at	dropoff_at
0	VT	2009-01-04 02:52:00	2009-01-04 03:02:00
1	VT	2009-01-04 03:31:00	2009-01-04 03:38:00
2	VT	2009-01-03 15:43:00	2009-01-03 15:57:00
3	DD	2009-01-01 20:52:58	2009-01-01 21:14:00
4	DD	2009-01-24 16:18:23	2009-01-24 16:24:56
5	DD	2009-01-16 22:35:59	2009-01-16 22:43:35
6	DD	2009-01-21 08:55:57	2009-01-21 09:05:42
7	VT	2009-01-04 04:31:00	2009-01-04 04:36:00
8	CM	2009-01-05 16:29:02	2009-01-05 16:40:21
9	CM	2009-01-05 18:53:13	2009-01-05 18:57:45



Remote SSH



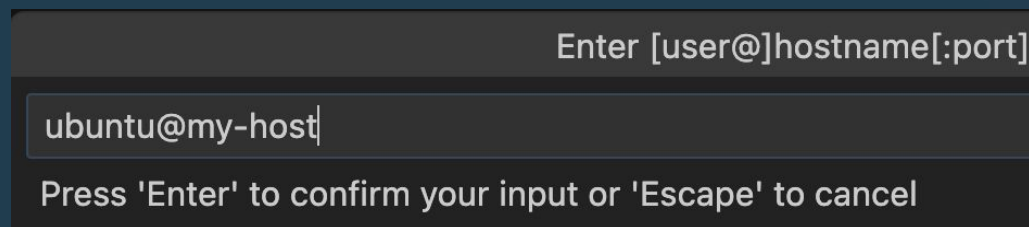
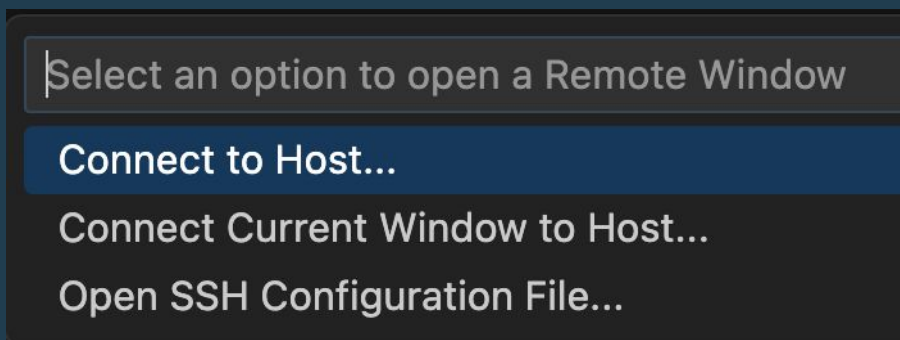
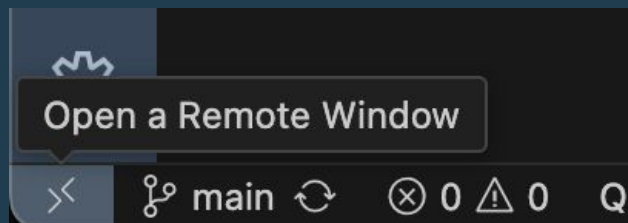
Who is this for?

- People who have access to another machine
 - Anything from a spare box to a high-performance cluster
- People who can set one up and maintain it
- Scrappy data science teams



REMOTE SSH

Connecting



analysis.py — ssh-demo [SSH: i-099684c376a0c179d]

Python 3.13.7 (Uv: ssh-demo) ssh-demo

EXPLORER

SSH-DEMO [SSH: i-099684c376a0c179d]

- .venv
- data
- .gitignore
- analysis.py
- download.py
- pyproject.toml
- README.md
- uv.lock

analysis.py

```
1 # Data from https://www.nyc.gov/site/tlc/about/tlc-trip-record-data.page
2 # hosted by Ursa Labs (see https://ursalabs.org/arrow-r-nightly/articles/dataset.htm
3
4 import ibis
5 from ibis import _
6 from plotnine import ggplot, aes, geom_bar, scale_x_continuous, ggtitle
7
8 # Set up any backend
9 con = ibis.duckdb.connect()
10
11 # Load an ibis table (this doesn't actually read the data yet)
12 t = con.read_parquet(
13     "data/**",
14     table_name="nyc taxi",
15     union_by_name=True,
16 )
17
18 # Inspect the table structure
19 t
20 len(t.columns)
21 t.count().execute()
22
23 # Find the mean tip percentage by number of passengers, among trips costing more than
24 tip_analysis = (
25     t.filter(
26         [
```

CONSOLE

~/.ssh-demo

- ... data=tip_analysis,
- ... mapping=aes(x="passenger_count", y="mean_tip_pct", fill="log_count"),
- ...)
- ... + geom_bar(stat="identity")
- ... + scale_x_continuous(breaks=range(15))
- ... + ggtitle("Tip % by passenger count (for expensive trips)")
- ...)
- ...
- 100%
- >>> |

SESSION CONNECTIONS HELP VIEWER

VARIABLES

Python 3.13.7 (Uv: ssh-demo) filter

VALUES

> _	-	Deferred
> con	<ibis.backends.duckdb.Backend objec...	Backend
> t	ibis.expr.types.relations.Table	ibis.Expr
> tip_analysis	ibis.expr.types.relations.Table	ibis.Expr

CLASSES

> aes	<class 'plotnine.mapping.aes.aes'>	type
> geom_bar	<class 'plotnine.geoms.geom_bar.geo...	Register

PLOTS

Tip % by passenger count (for expensive trips)

passenger_count	mean_tip_pct	log_count
1	13.0	12
2	12.5	11
3	12.0	10
4	11.0	9
5	14.0	8
6	14.5	7
7	15.0	6
8	14.0	6
9	12.5	6

Ln 49, Col 1 Spaces: 4 UTF-8 LF {} Python



NewOpen

analysis.py — ssh-demo [SSH: i-099684c376a0c179d]

Python 3.13.7 (Uv: ssh-demo)ssh-demo

EXPLORER

SSH-DEMO [SSH: i-099684c376a0c179d]

- .venv
- data
- .gitignore
- analysis.py
- download.py
- pyproject.toml
- README.md
- uv.lock

analysis.py

```
1 # Data from https://www.nyc.gov/site/tlc/about/tlc-trip-record-data.page
2 # hosted by Ursa Labs (see https://ursalabs.org/arrow-r-nightly/articles/dataset.htm
3
4 import ibis
5 from ibis import _
6 from plotnine import ggplot, aes, geom_bar, scale_x_continuous, ggtitle
7
8 # Set up any backend
9 con = ibis.duckdb.connect()
10
11 # Load an ibis table (this doesn't actually read the data yet)
12 t = con.read_parquet(
13     "data/**",
14     table_name="nyc taxi",
15     union_by_name=True,
16 )
17
18 # Inspect the table structure
19 t
20 len(t.columns)
21 t.count().execute()
22
23 # Find the mean tip percentage by number of passengers, among trips costing more than
24 tip_analysis = (
25     t.filter(
26         [
```

SESSIONCONNECTIONSHELPERVIEWER

VARIABLES

Python 3.13.7 (Uv: ssh-demo)filter

VALUES

> _	-	Deferred
> con	<ibis.backends.duckdb.Backend objec...	Backend
> t	ibis.expr.types.relations.Table	ibis.Expr
> tip_analysis	ibis.expr.types.relations.Table	ibis.Expr

CLASSES

> aes	<class 'plotnine.mapping.aes.aes'>	type
> geom_bar	<class 'plotnine.geoms.geom_bar.geo...	Register

PLOTS

Tip % by passenger count (for expensive trips)

CONSOLETERMINALPROBLEMSOUTPUTPORTS2DEBUG CONSOLE

~/ssh-demo

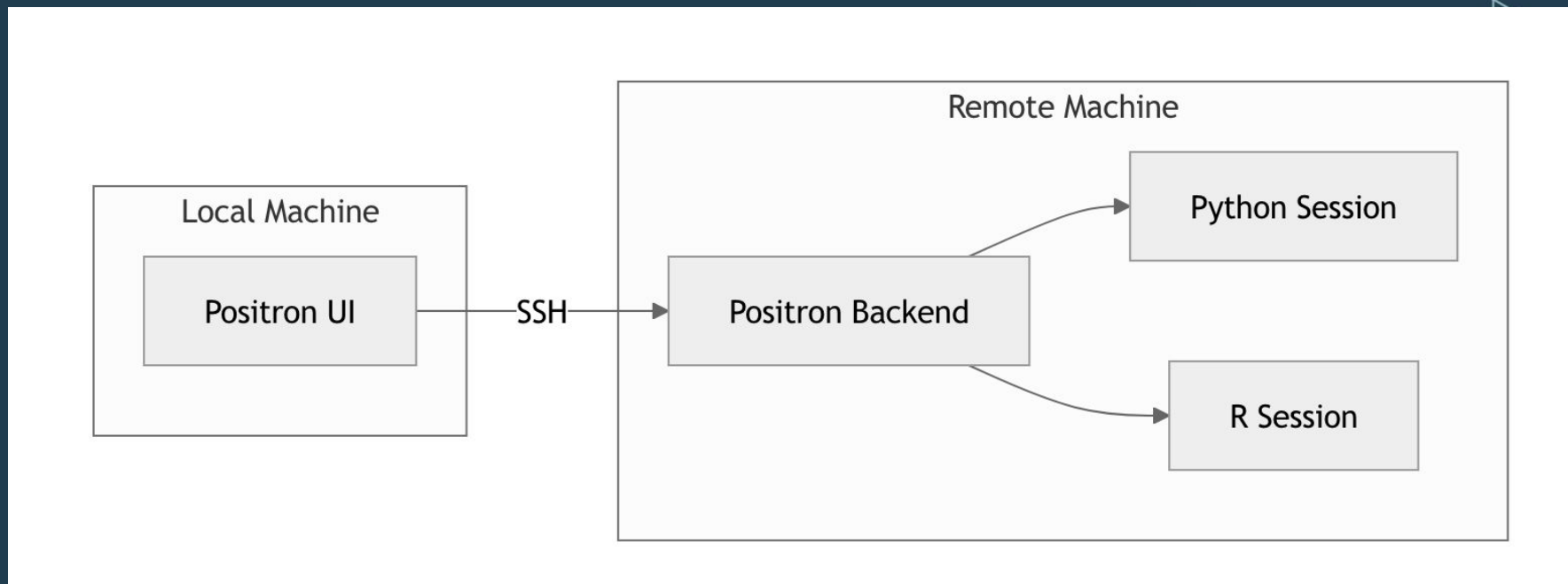
```
... data=tip_analysis,
... mapping=aes(x="passenger_count", y="mean_tip_pct", fill="log_count"),
... )
... + geom_bar(stat="identity")
... + scale_x_continuous(breaks=range(15))
... + ggtitle("Tip % by passenger count (for expensive trips)")
... )
...
100%
>>> |
```

OUTLINE

TIMELINE

posi
conf (2025)

Architecture

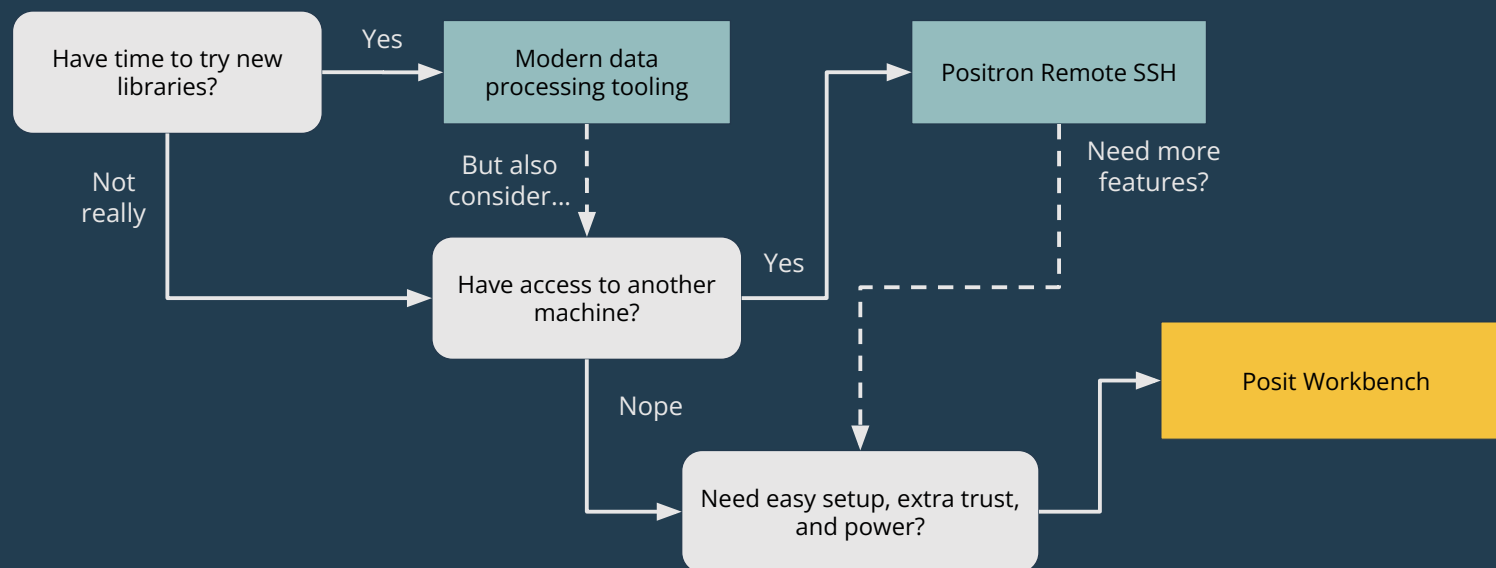


Benefits

- It's free to use, including at work
- You don't have to change any code
 - Use pandas on a huge machine if you really want!
- You can disconnect and reconnect while running code

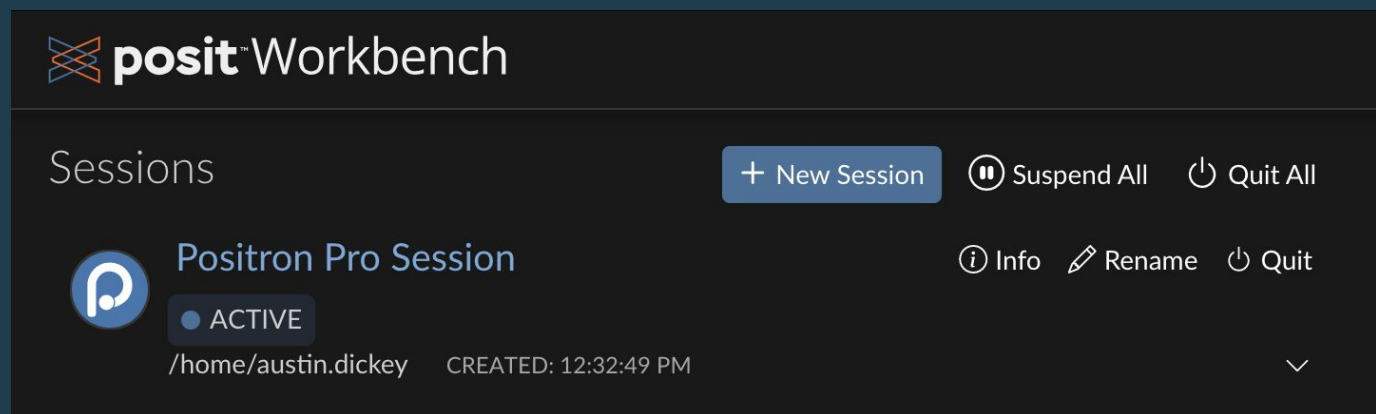


Posit Workbench



Who is this for?

- People who need a lot more power and scaling
- Industries with high compliance and trust standards
- Teams looking for standardized environments and data access



Benefits

- Like SSH and more
- Cluster management, security/auditing, consistent development environments
- Working with partners like Snowflake



Want to learn more?

- "Posit's Commercial Products" - today at 1:00pm
- "Enterprise Data Platforms" - tomorrow at 2:40pm
- I'll hang out in the Positron Lounge tomorrow morning

