# Winning Strategies:

# Powering a World Series Victory with Airflow Orchestration
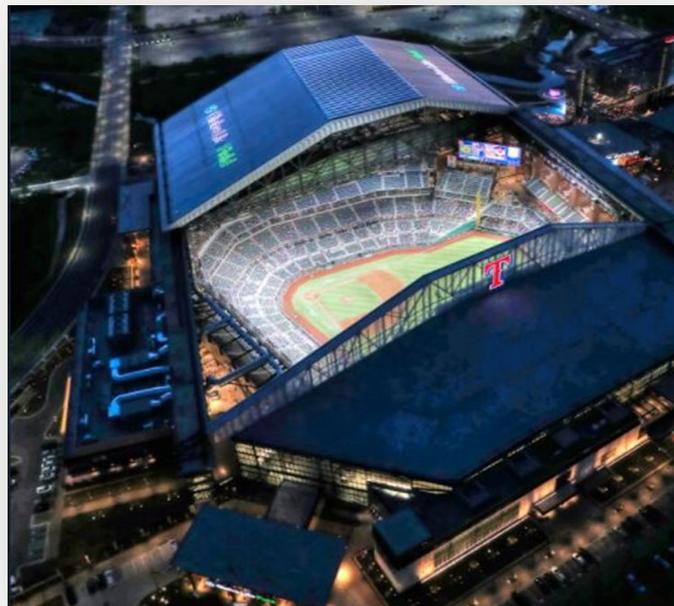
Oliver Dykstra

# Who Am I

## Oliver Dykstra

Data Engineer, R&D
Texas Rangers Baseball Club
Joined the club in 2022
odykstra@texasrangers.com

# What We're Talking About

- History of Baseball Statistics

- Big Data Baseball

- The Winning Formula:

    - Upgrading with Airflow

    - Astronomer Take the Wheel

- Putting it all together:

    ***World Series Victory***

# A Brief History of Baseball Statistics

That's all baseball is, is numbers; it's run by numbers, averages, percentages and odds...
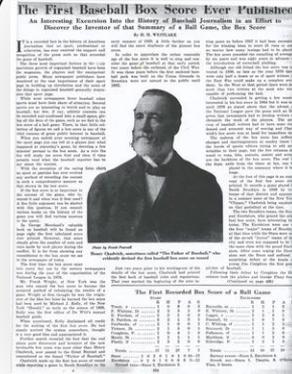
– Rollie Fingers
(Hall of Fame Pitcher)

# Think Outside the Box(score)

## A Brief History of Baseball Statistics

**1859** - Henry Chadwick publishes the first Box Score, a set of statistics compiling the runs, hits, outs, assists and errors.

**~~ 100 years later ~~**

**1952** - Topps adds full statistics lines on the back of their annual baseball cards.
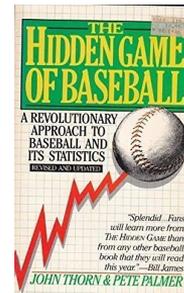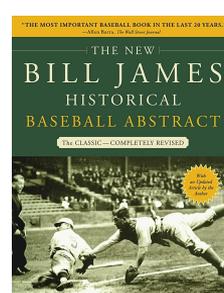
# Think Outside the Box(score)

## A Brief History of Baseball Statistics

**1977** - Bill James publishes his first "Baseball Abstract" which becomes a national bestseller in the early 80s.

**1989** - Retrosheet begins massive compilation and online publishing of old box scores and play-by-plays, allowing droves of historical research never before possible.

**1996** - Baseball Prospectus begins publication of their annual and website - introduces statistics community to VORP, PECOTA, Pitcher Abuse points and more.
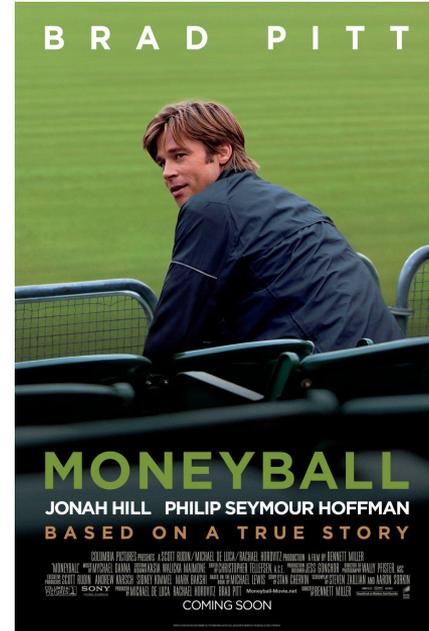
# Think Outside the Box(score)

## A Brief History of Baseball Statistics

**2001 - 2003** - Billy Beane's Moneyball A's use data-driven insights to identify market inefficiencies within the game of baseball.

**2003** - Michael Lewis publishes his book, "Moneyball: The Art of Winning an Unfair Game" about those A's.

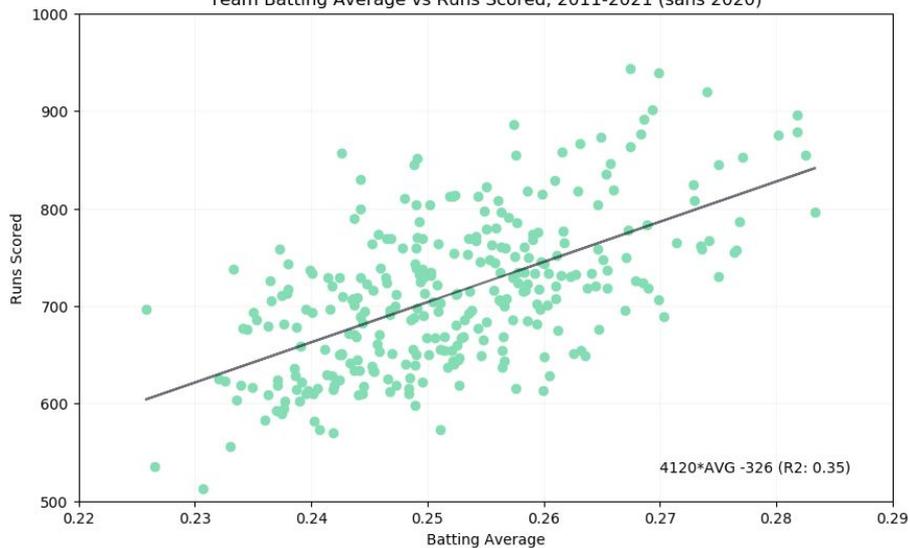**2011** - Brad Pitt stars in the movie adaptation of the same name.
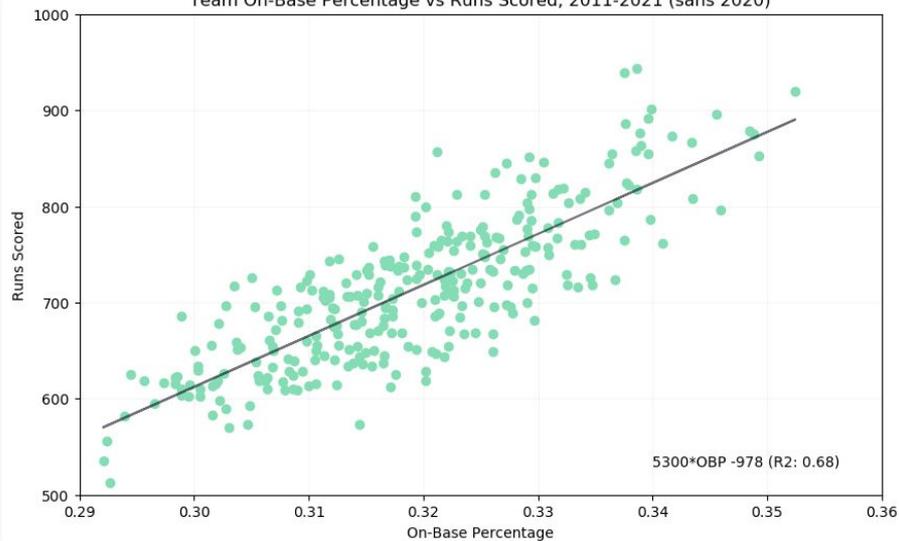
# Think Outside the Box(score)
## A Brief History of Baseball Statistics



Team Batting Average vs Runs Scored, 2011-2021 (sans 2020)

4120*AVG -326 (R2: 0.35)



Team On-Base Percentage vs Runs Scored, 2011-2021 (sans 2020)

5300*OBP -978 (R2: 0.68)

# Think Outside the Box(score)

## A Brief History of Baseball Statistics

Billy Beane identified a **market inefficiency**.

The market priced players with high batting averages higher than those with high on-base percentages. However, on-base percentage has a **higher correlation** to total runs scored.

The Oakland A's used this information to acquire players undervalued by the market that could help them compete with higher payroll teams.

This data-driven decision **disrupted the industry** and left a legacy far beyond baseball.

Big Data Baseball

"If you challenge conventional wisdom, you will find ways to do things much better than they are currently done."

–  Bill James

(Founder of Sabermetrics)

# The Statcast Revolution

## Big Data Baseball

**2001–2002**: Moneyball, Billy Beane, Oakland A's identify data-driven market inefficiencies.

**2015:** Statcast Debut. Radar + HD Video measures all action on the field, per pitch.

**2020:** Statcast switches from TrackMan to Hawk–Eye as its technology provider.

**2022:** Statcast deployed to AAA. Widespread MiLB adoption planned.

**2006:** PITCHf/x Ball Tracking Debut. Spin rates, velocity, and movement all tracked.

**2017:** Statcast switches from PITCHf/x to TrackMan as its technology provider.

**2021:** Pose Tracking and FieldVision debut. Skeleton and body movements tracked.

# The Statcast Revolution

## Big Data Baseball

# The Statcast Revolution

## Big Data Baseball

# The Statcast Revolution

## Big Data Baseball

# The Statcast Revolution

**Big Data Baseball**

# The Statcast Revolution

## Big Data Baseball

# The Statcast Revolution

## Big Data Baseball

# The Statcast Revolution

**Big Data Baseball**

# The Statcast Revolution

## Big Data Baseball

# The Winning Formula: Upgrading With Airflow

"It's math-y, but there's still the whole arts-and-science debate [with big data]. I'd argue there is an art to that sort of stuff."

Mike Fitzgerald
(VP of R&D for the Diamondbacks)

# The Old Way

## The Winning Formula

Our first data solution failed to adapt as we started ingesting big data.

Problems:
- Cron Scheduled
- Hard to troubleshoot
- Tough to explain

# The Potential

## The Winning Formula

- Consumers want data fast
- Individualized requirements
- Global consumers and schedules

# New Data Source Ingestion

## The Winning Formula

Statcast **bat tracking** is available beginning with the 2024 season.

Since different parts of a bat can move at different speeds, an individual swing's speed is measured at the point six inches from the head of the bat, what is popularly called *"the sweet-spot."*

**Swing length** tracks the sum distance traveled by the head of the bat in XYZ space from the start of data until contact point.



Bat Speed Distribution Comparison

C. Seager (2024 - Left)
M. Semien (2024 - Right)
R. Acuña Jr. (2024 - Right)
League (2024)

Swings
Bat Speed

LOCK SCALES | YES ▼    SCALE | % TOTAL SWINGS ▼    SWING % ▼

FIND A PLAYER BY LAST NAME...

# A Unified Data Platform
## The Winning Formula

# Upgrade With Airflow

## The Winning Formula

# Custom Monitoring
The Winning Formula

# Named Dynamic Tasks v2.9

## The Winning Formula



```python
with DAG('GAME_Events',
        start_date=datetime.datetime(2020, 1, 1),
        schedule="30 6 * * *",
        default_args=default_args,
        catchup=False,
        max_active_runs=1,
        tags=["event", "pitches"]
        ) as dag:

    ingest_game_ids = PythonOperator(
        task_id="ingest_game_ids",
        python_callable=synergy.get_game_ids
    )

    game_ids_to_delta = DatabricksRunNowOperator(
        task_id='game_ids_to_delta',
        databricks_conn_id='databricks',
        job_id=903119014908886,
        notebook_params = {'file_path' : f"s3://{BUCKET}/"
    )

    get_events_for_game_id = PythonOperator.partial(
        task_id=f"get_events_by_game_id",
        python_callable=synergy.get_events_by_game_id,
        map_index_template="{{ task.op_args[0] }}"
    ).expand(op_args=ingest_game_ids.output)
```

# Data Driven Scheduling v2.9

## The Winning Formula

```python
current_time = datetime.now(tz=timezone.utc)

fact_pitch_dataset = Dataset("fact_pitch")
dim_video_dataset = Dataset("dim_video")
chipotle_dataset = Dataset("chipotle")
mlb_analytics_dataset = Dataset("mlb_analytics")

# trigger dataset if time is within a window
if current_time > datetime(current_time.year, current_time.month, current_time.day, 7, 00, tzinfo=timezone.utc) \
    and current_time < datetime(current_time.year, current_time.month, current_time.day, 10, 15, tzinfo=timezone.utc):

    whpm_dataset = Dataset("morning_whpm")
else:
    whpm_dataset = Dataset("whpm")

with DAG('WHPM_Merge_Pitches',
    start_date=datetime(2020, 7, 1),
    schedule=((fact_pitch_dataset & dim_video_dataset) | chipotle_dataset),
    default_args=default_args,
    catchup=False,
    max_active_runs=1,
    tags=["whpm"]
    ) as dag:
```

# Data Driven Scheduling v2.9

## The Winning Formula

```python
current_time = datetime.now(tz=timezone.utc)

fact_pitch_dataset = Dataset("fact_pitch")
dim_video_dataset = Dataset("dim_video")
chipotle_dataset = Dataset("chipotle")
mlb_analytics_dataset = Dataset("mlb_analytics")

# trigger dataset if time is within a window
if current_time > datetime(current_time.year, current_time.month, current_time.day, 7, 00, tzinfo=timezone.utc) \
    and current_time < datetime(current_time.year, current_time.month, current_time.day, 10, 15, tzinfo=timezone.utc):

    whpm_dataset = Dataset("morning_whpm")
else:
    whpm_dataset = Dataset("whpm")

with DAG('WHPM_Merge_Pitches',
    start_date=datetime(2020, 7, 1),
    schedule=((fact_pitch_dataset & dim_video_dataset) | chipotle_dataset),
    default_args=default_args,
    catchup=False,
    max_active_runs=1,
    tags=["whpm"],
    ) as dag:
```
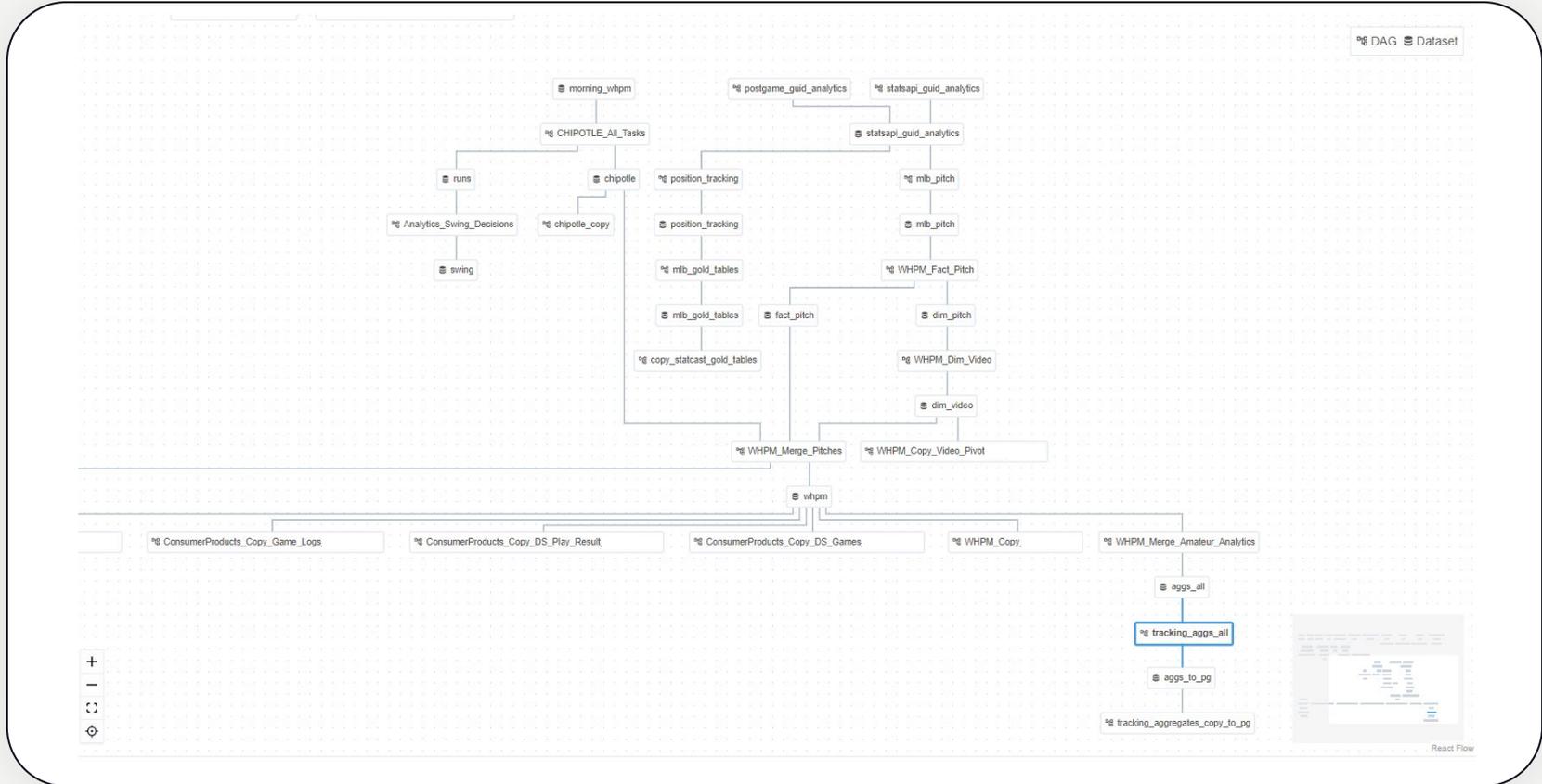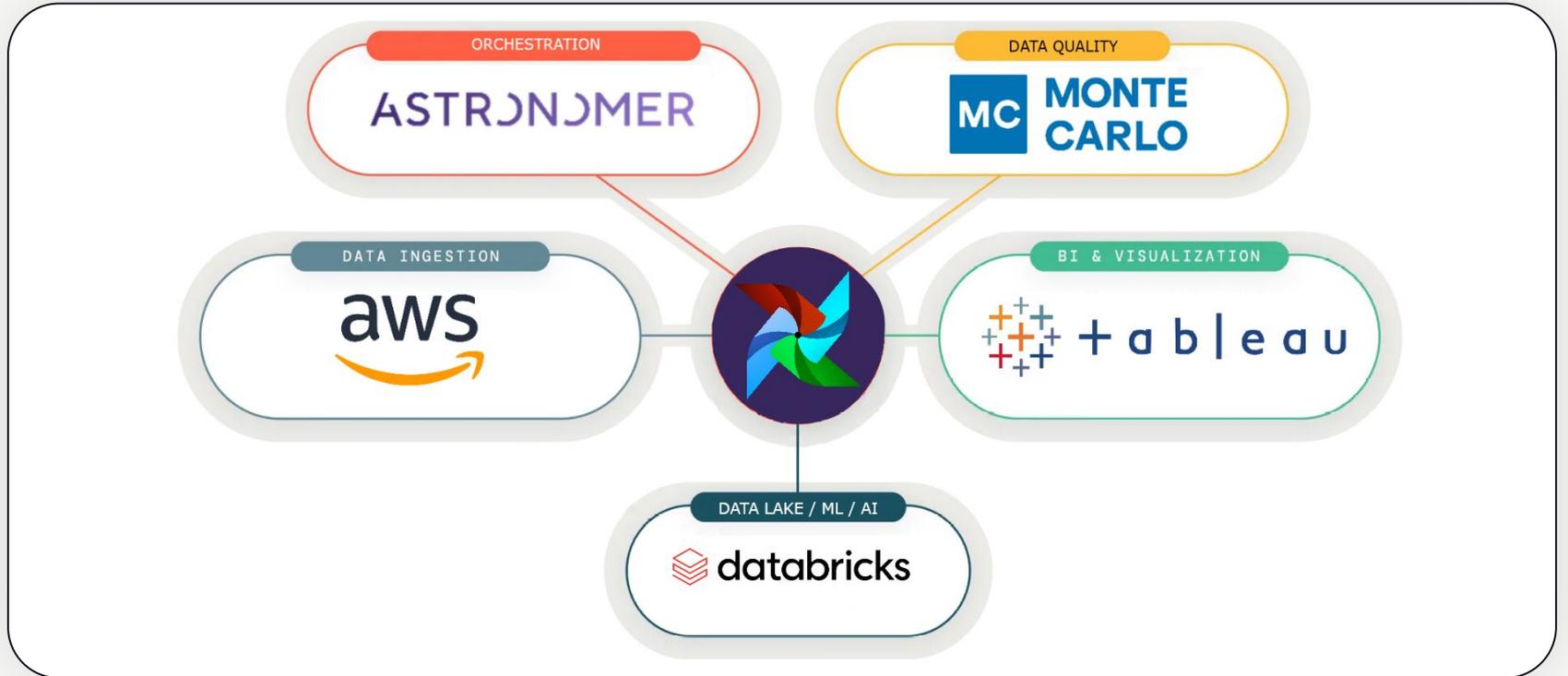
# Data Driven Scheduling v2.9

## The Winning Formula

```python
current_time = datetime.now(tz=timezone.utc)

fact_pitch_dataset = Dataset("fact_pitch")
dim_video_dataset = Dataset("dim_video")
chipotle_dataset = Dataset("chipotle")
mlb_analytics_dataset = Dataset("mlb_analytics")

# trigger dataset if time is within a window
if current_time > datetime(current_time.year, current_time.month, current_time.day, 7, 00, tzinfo=timezone.utc) \
    and current_time < datetime(current_time.year, current_time.month, current_time.day, 10, 15, tzinfo=timezone.utc):

    whpm_dataset = Dataset("morning_whpm")
else:
    whpm_dataset = Dataset("whpm")

with DAG('WHPM_Merge_Pitches',
    start_date=datetime(2020, 7, 1),
    schedule=((fact_pitch_dataset & dim_video_dataset) | chipotle_dataset),
    default_args=default_args,
    catchup=False,
    max_active_runs=1,
    tags=["whpm"],
    ) as dag:
```

# Data Driven Scheduling v2.9

## The Winning Formula

```python
current_time = datetime.now(tz=timezone.utc)

fact_pitch_dataset = Dataset("fact_pitch")
dim_video_dataset = Dataset("dim_video")
chipotle_dataset = Dataset("chipotle")
mlb_analytics_dataset = Dataset("mlb_analytics")

# trigger dataset if time is within a window
if current_time > datetime(current_time.year, current_time.month, current_time.day, 7, 00, tzinfo=timezone.utc) \
    and current_time < datetime(current_time.year, current_time.month, current_time.day, 10, 15, tzinfo=timezone.utc):

    whpm_dataset = Dataset("morning_whpm")
else:
    whpm_dataset = Dataset("whpm")

with DAG('WHPM_Merge_Pitches',
    start_date=datetime(2020, 7, 1),
    schedule=((fact_pitch_dataset & dim_video_dataset) | chipotle_dataset),
    default_args=default_args,
    catchup=False,
    max_active_runs=1,
    tags=["whpm"],
    ) as dag:
```

# Data Driven Scheduling v2.9

## The Winning Formula

# Partner Solutions

## The Winning Formula

# Open Source

## The Winning Formula

"Half of it's art... it's creativity, and then half of it is just knowing the data you're working with and being able to manipulate it in the direction that will benefit the player."

Brian Bannister
(Director of Pitching, Chicago White Sox)

# Managed Airflow Environment
## The Winning Formula

# Managed Airflow Environment
## The Winning Formula

# Managed Airflow Environment

## The Winning Formula

# Managed Airflow Environment

## The Winning Formula

# Managed Airflow Environment

## The Winning Formula

CI/CD with Astronomer–
- Guaranteed Code Review
- Automatic Tests and Deployment
- Easy Dev Promotion Process

# The Statcast Revolution

## The Winning Formula
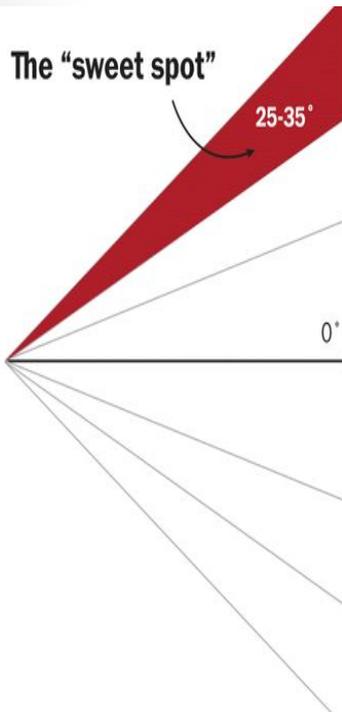
# Putting It All Together: World Series Victory

"If you can make yourself half a percent better... then that's a win. It can be the difference between making the playoffs and not making the playoffs."
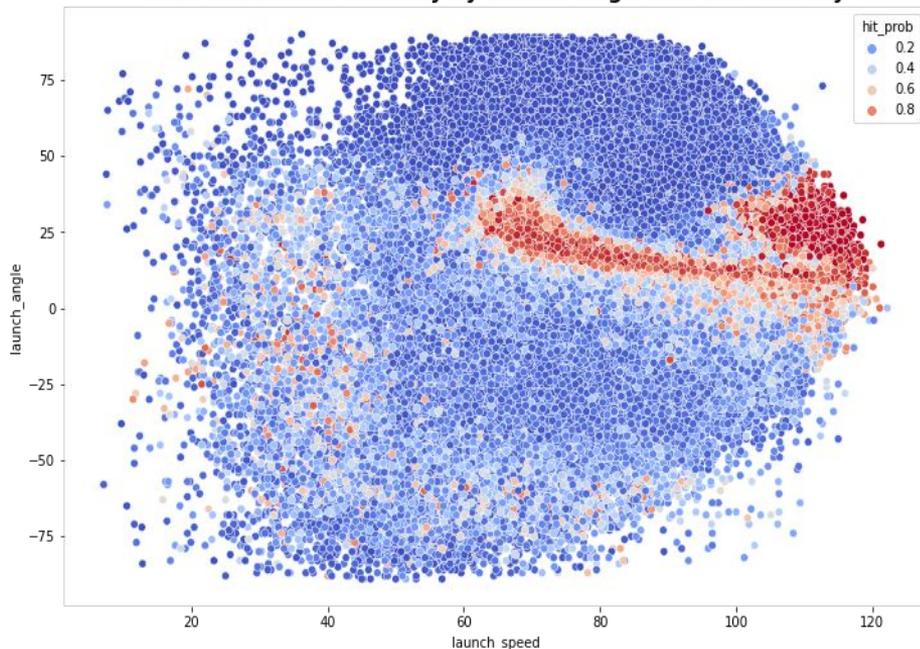
Ryan Murray
(Sr. Director of Baseball R&D,
Texas Rangers)

# The Launch Angle Revolution
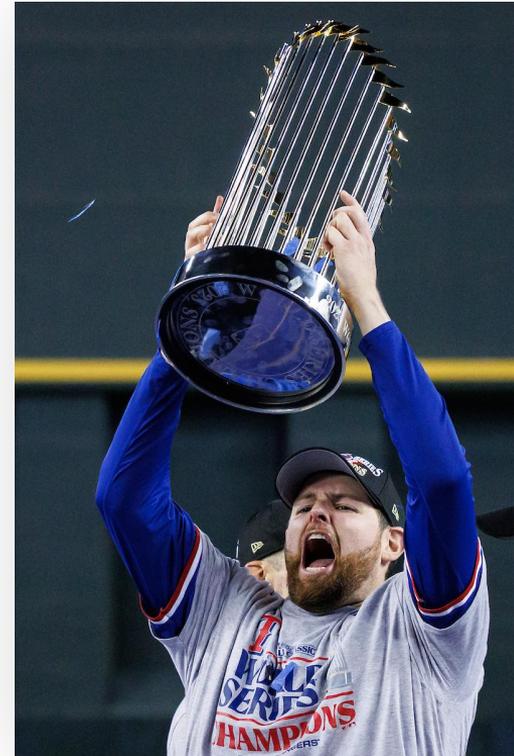
## Putting It All Together



The "sweet spot" 25-35°  0°

2019-2021 Hit Probability by Launch Angle and Exit Velocity

# Reports and Visualizations

## Putting It All Together

# World Series Victory

## Putting It All Together

# World Series Victory

## Putting It All Together

# Questions?