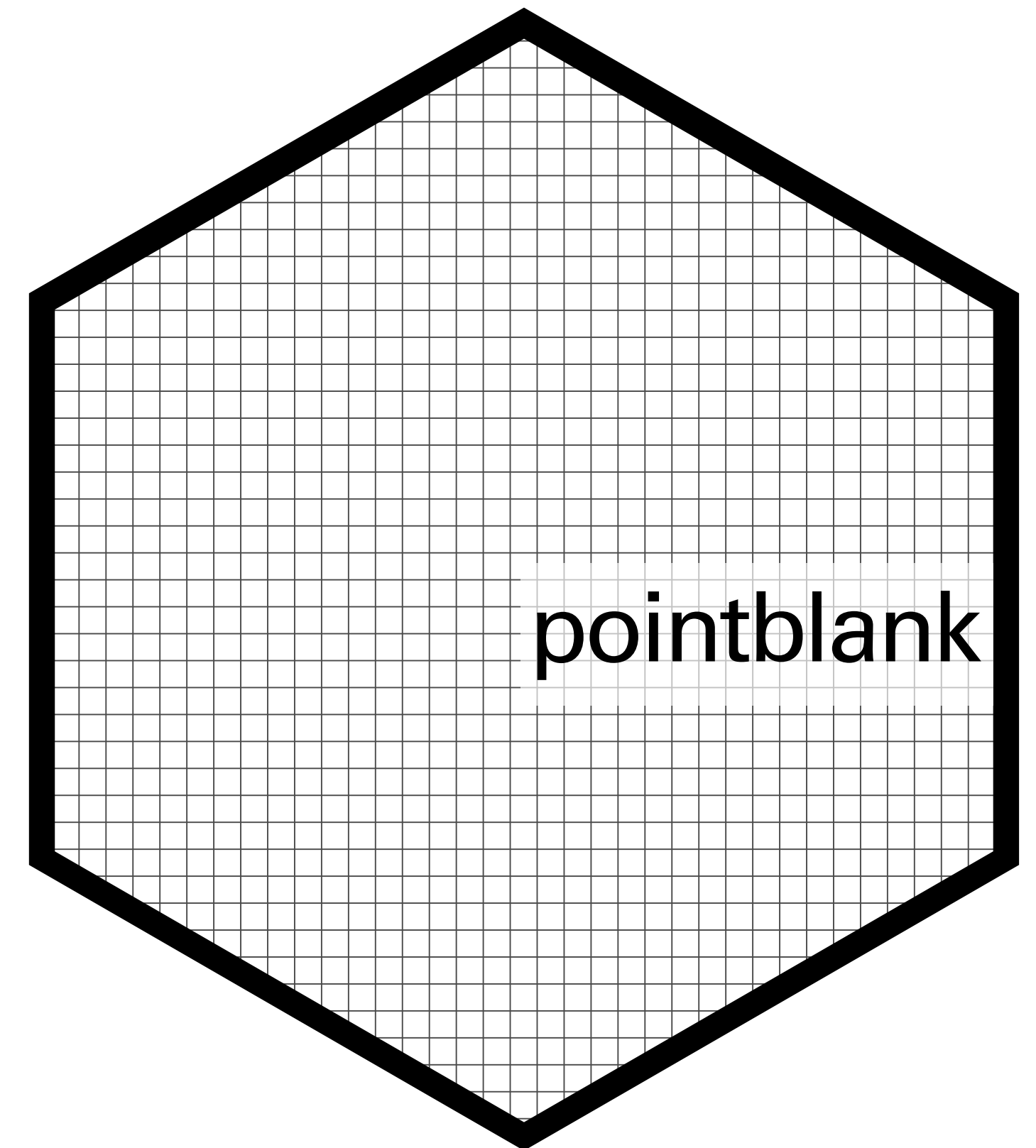
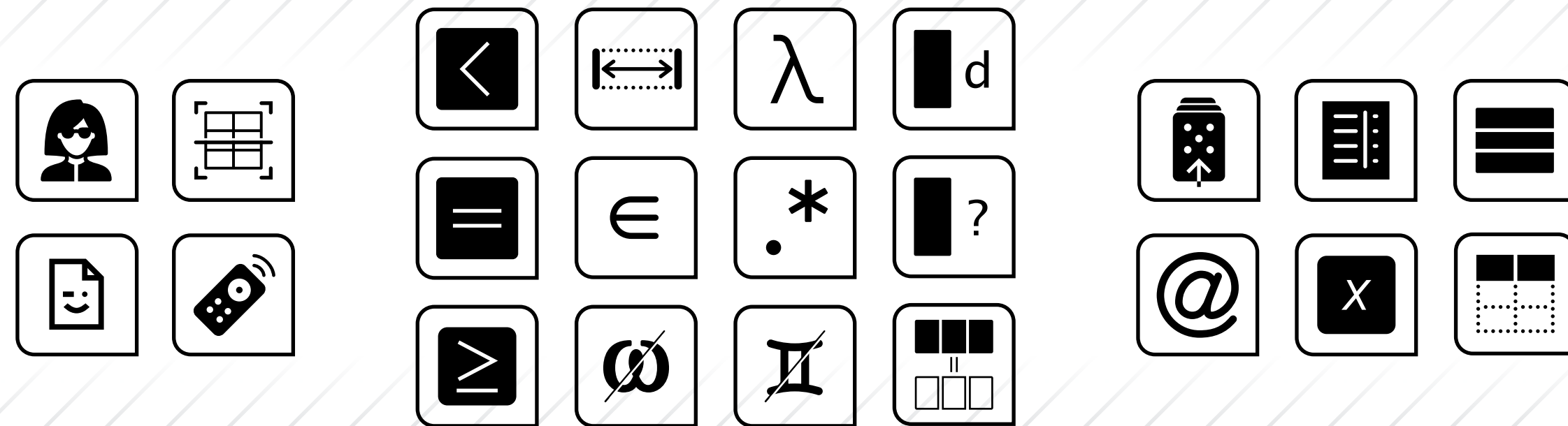


The **pointblank** R Package



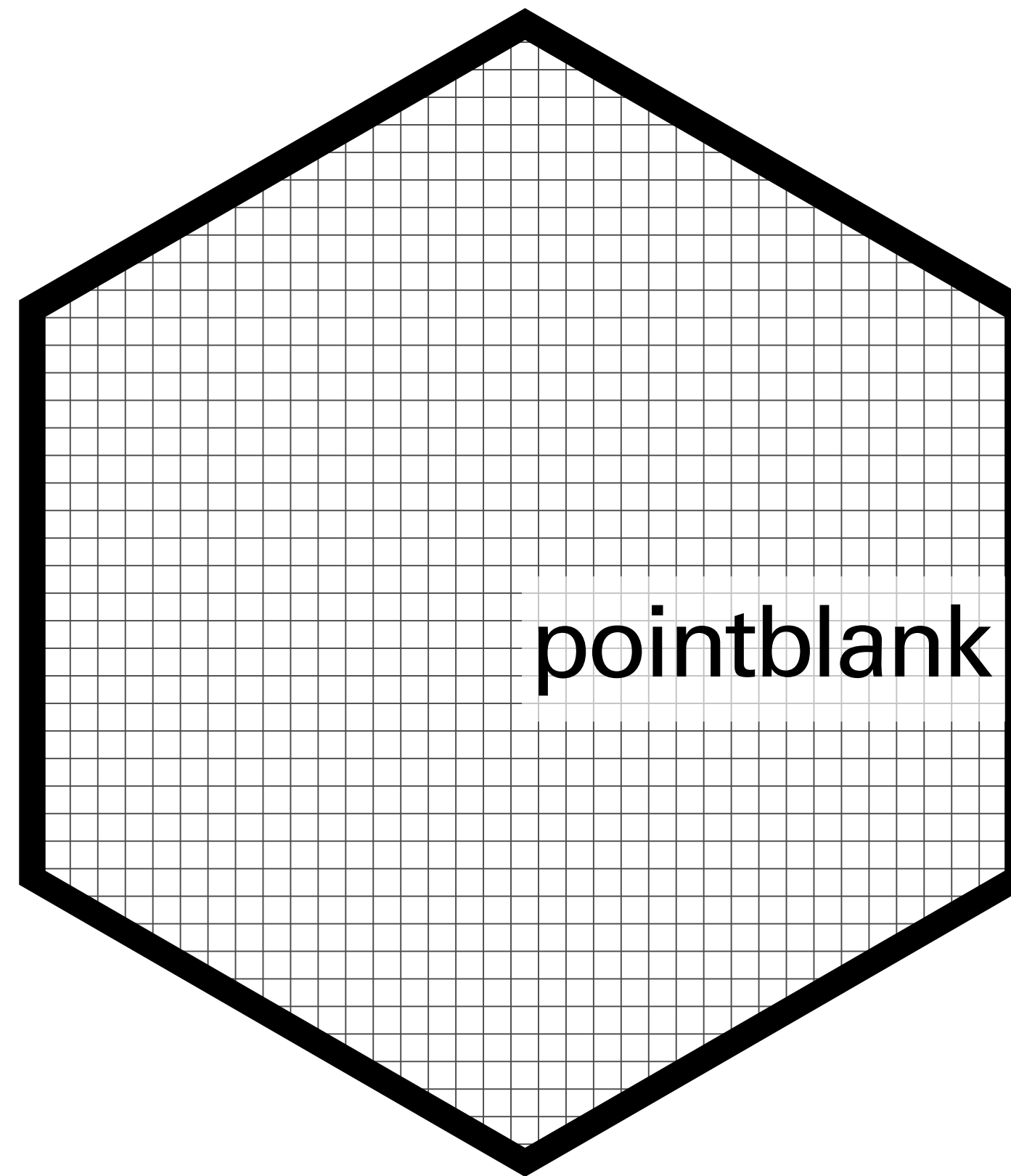
rich-iannone



@riannone



rich@rstudio.com



It's a package for validating data tables.

Because, more often than not, they are not free of errors.

Some Design Goals of **pointblank**

Make it work with all sorts of different tables, like **local tables** and **database tables**.



data.frame tbl_dbi*
tbl_df tbl_spark

*tested with: **MySQL**, **SQLite**, and **PostgreSQL**.

Provide a way to **understand** new datasets.

Overview of **dplyr::storms**

Overview		Reproducibility	
Table Overview		Column Types	
Columns	13	numeric	7
Rows	10010	integer	3
NA s	13056 (10%)	character	2
Duplicate Rows	0	ordered	1

Variables

name character	Distinct	198 (2%)	
	NA s	0	
	Inf / -Inf	0	
Toggle details			

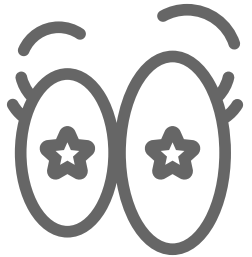
year numeric	Distinct	41 (0.4%)	Mean	1998.11
	NA s	0	Minimum	1975
	Inf / -Inf	0	Maximum	2015
Toggle details				

Create a toolset that can accommodate a lot of different **data validation workflows**.

- 1 DATA QUALITY
- 2 ETL or ANALYSIS SCRIPT
- 3 UNIT TESTING
- 4 CONDITIONAL CODE
- 5 TABLE SCAN
- 6 RMarkdown VALIDATION

THERE'S MORE

Keep trying to **make it easy** to use the package, with clear docs and examples.



Have reporting outputs be really nice to look at and **useful to everyone** in an organization.



Have reporting outputs translated to multiple **spoken languages**.
EN ■ FR ■ DE ■ IT ■ ES

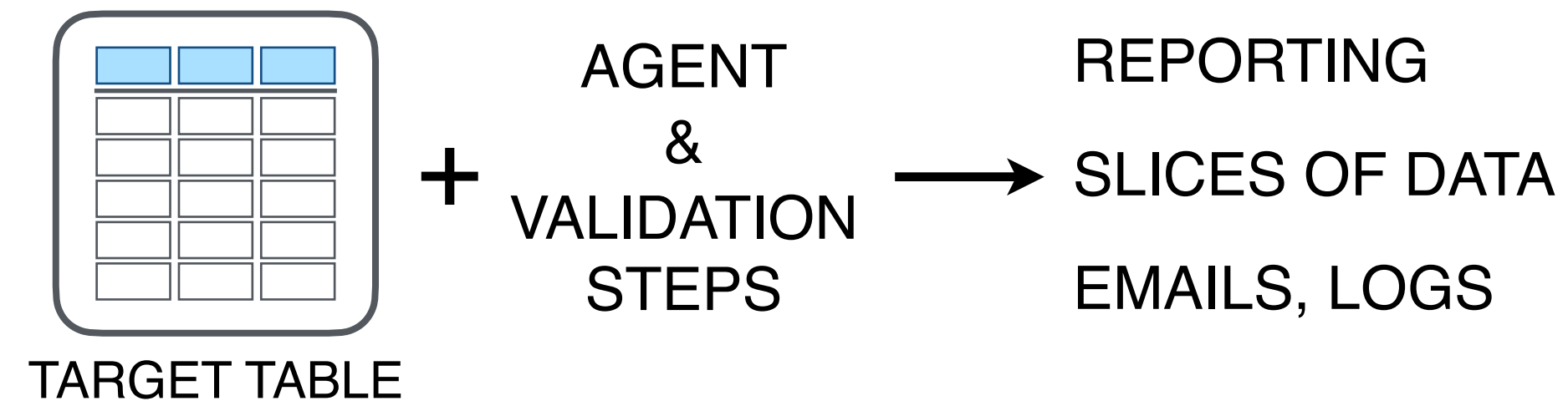


→ and more to come!

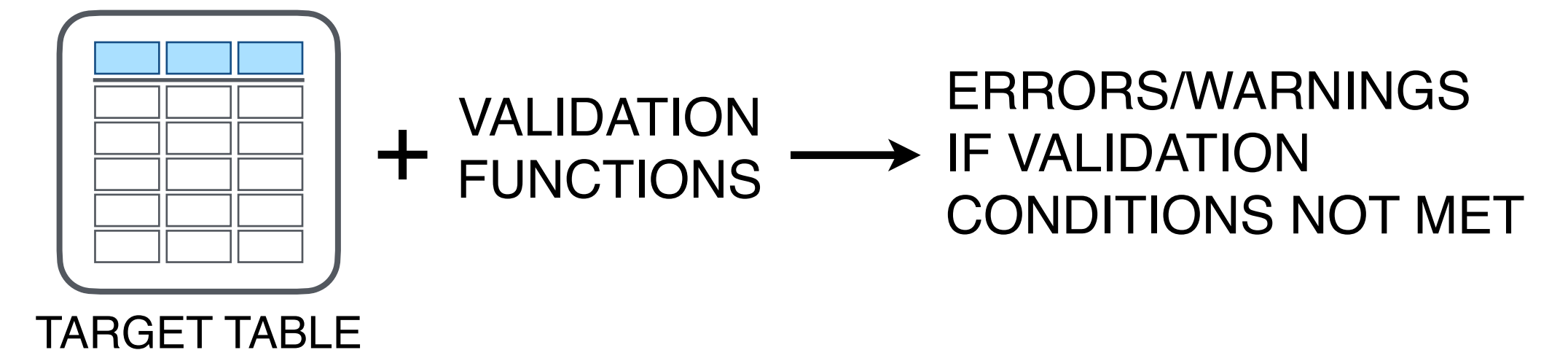
Useful Workflows in **pointblank**

MAIN WORKFLOWS

① DATA QUALITY



② ETL or ANALYSIS SCRIPT



SECONDARY WORKFLOWS

③ UNIT TESTING

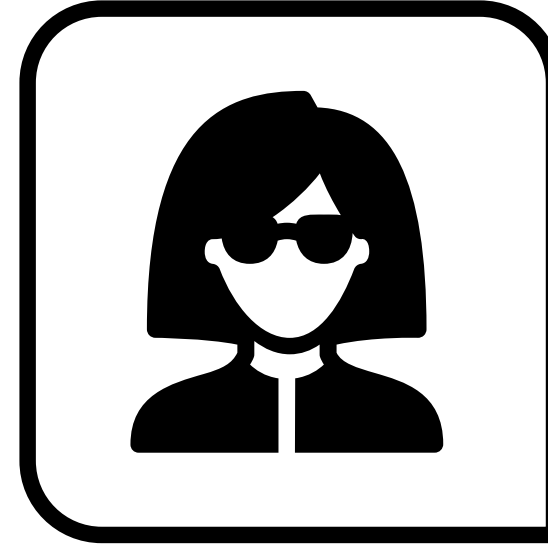
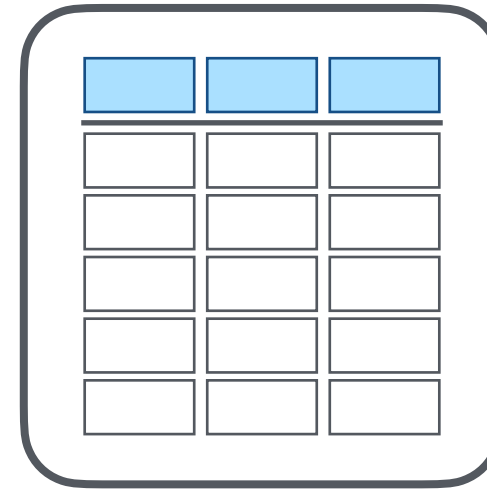
You can validate your data tables in your unit tests. It's just like **testthat** actually. *Great* for packages.

④ CONDITIONAL CODE

Get **TRUE** or **FALSE** based on a data validation. This can be useful **in R code**.

The **pointblank** Data Quality Workflow^①

The **agent** is given
the **target table**...



`create_agent()`

The **agent** is an
integral part of the
data quality workflow.

The **pointblank** Data Quality Workflow^①

The **agent** is given
the **target table**...

actions

end_fns

lang

...and some directives
on interrogation.



create_agent()

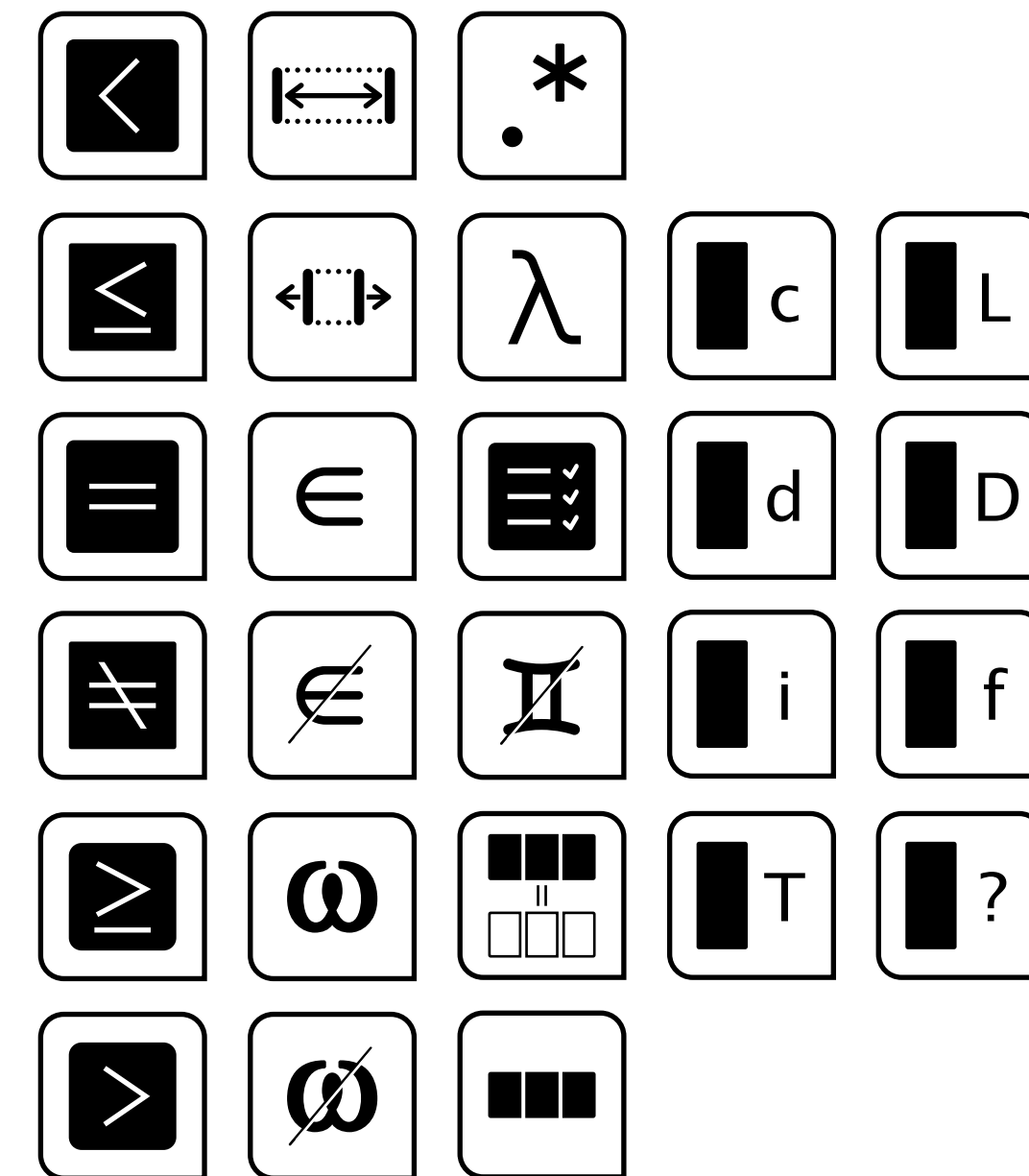
The **agent** is an
integral part of the
data quality workflow.

The **pointblank** Data Quality Workflow^①



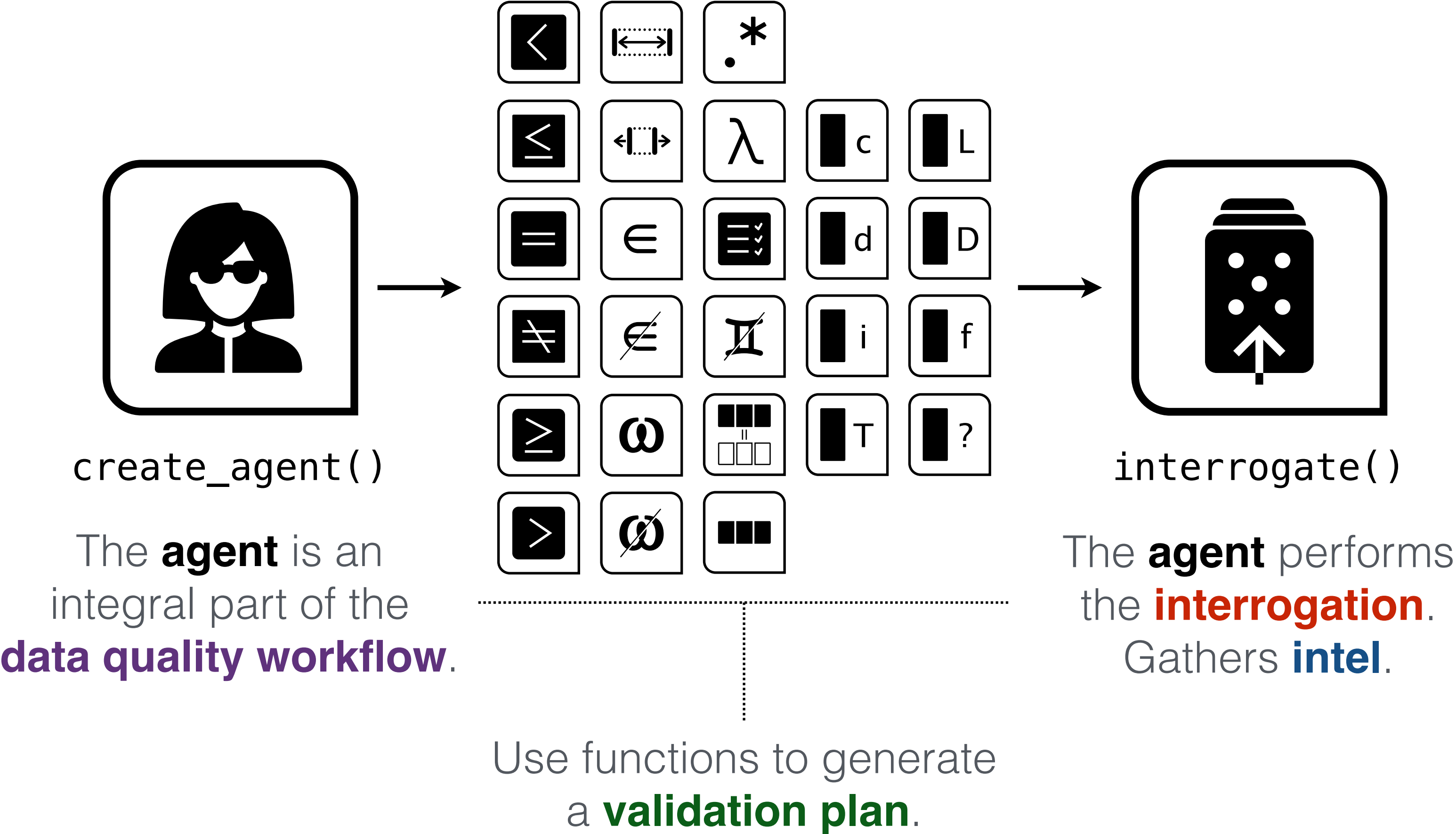
`create_agent()`

The **agent** is an integral part of the **data quality workflow**.

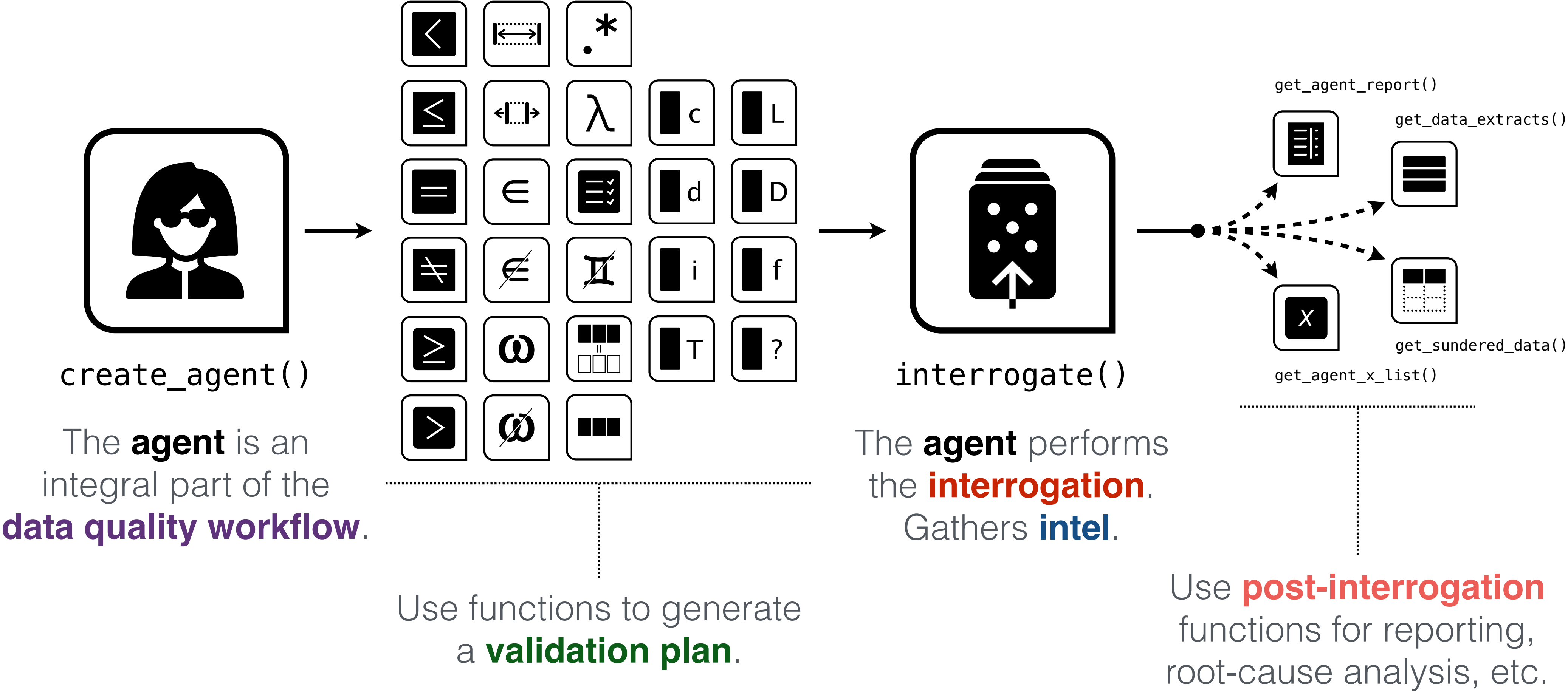


Use functions to generate a **validation plan**.

The pointblank Data Quality Workflow^①

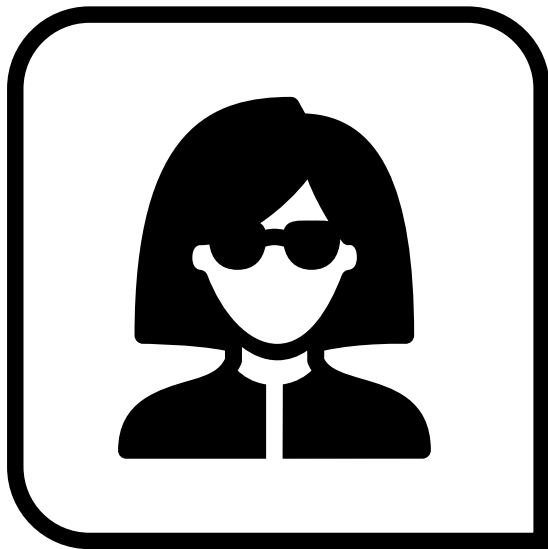
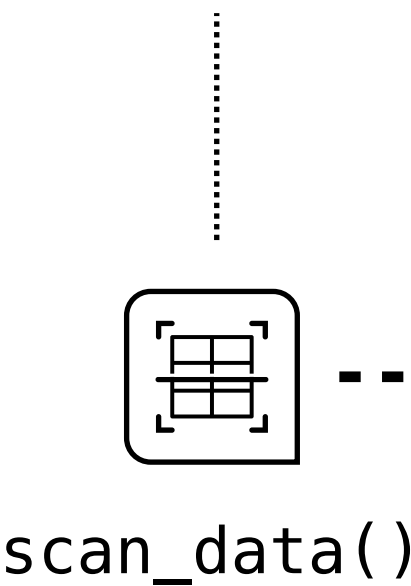


The pointblank Data Quality Workflow^①



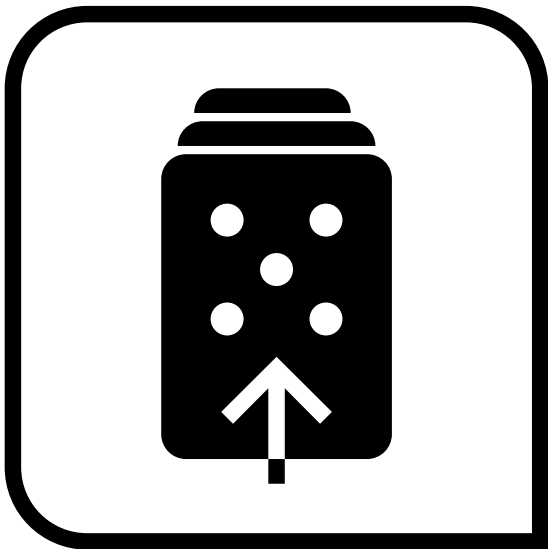
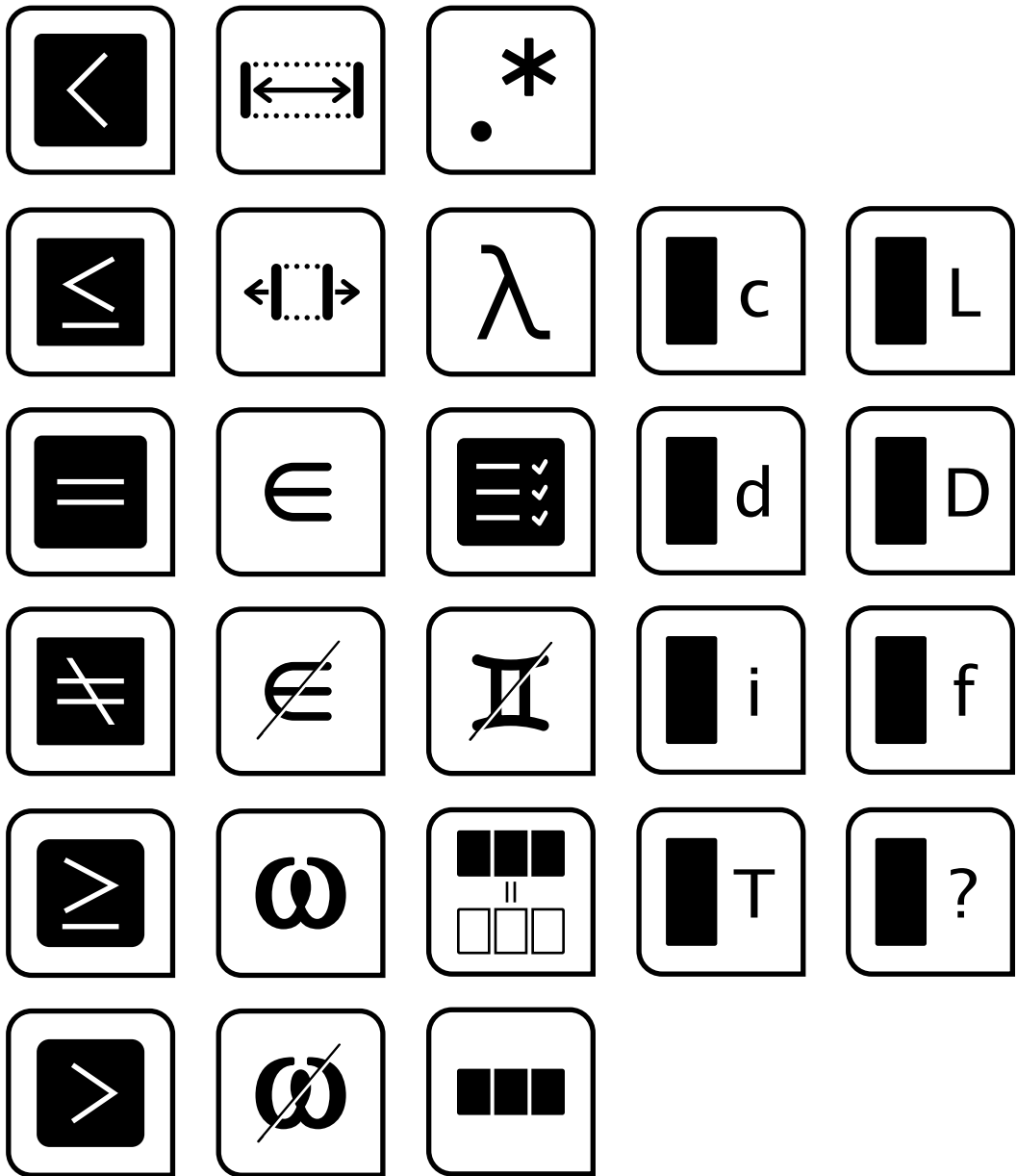
The **pointblank** Data Quality Workflow^①

Understand your data more with a **table scan**.



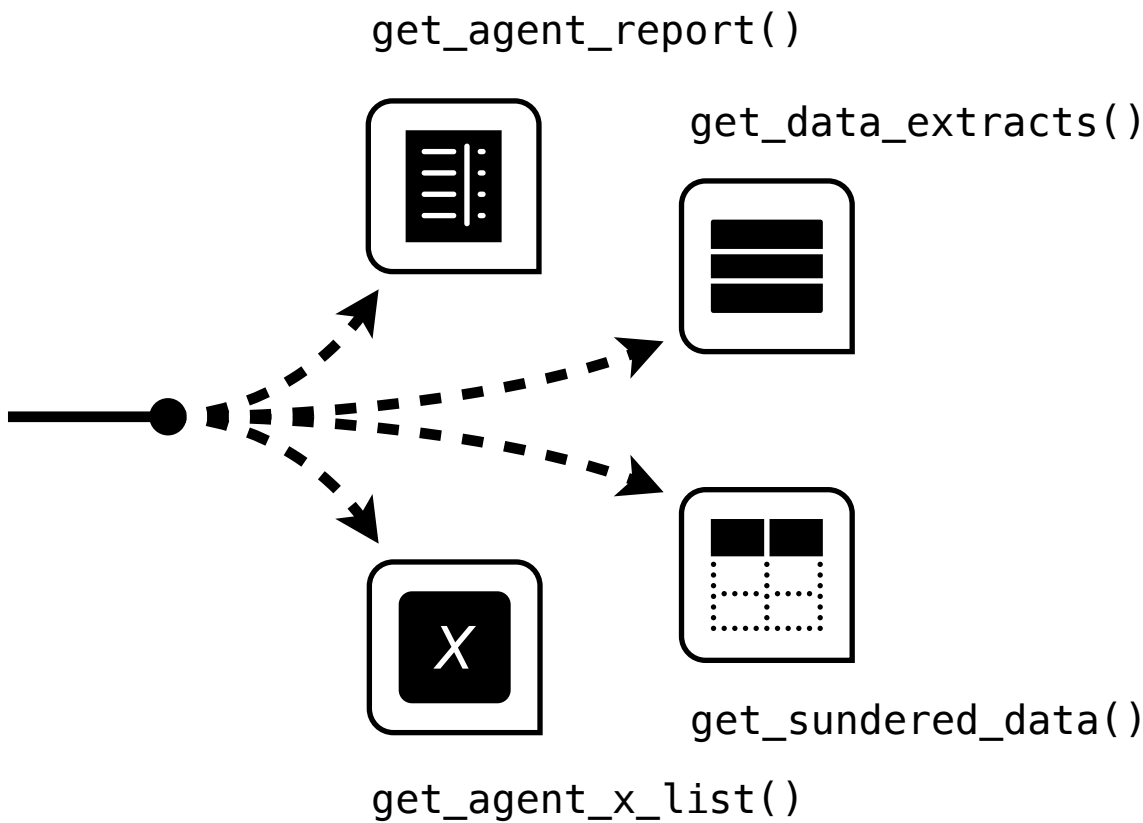
create_agent()

The **agent** is an integral part of the **data quality workflow**.



interrogate()

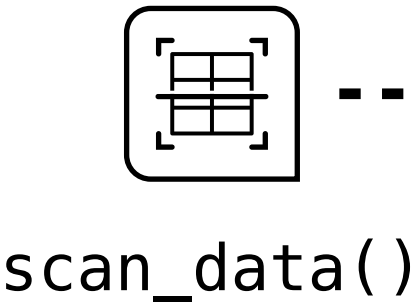
The **agent** performs the **interrogation**. Gathers **intel**.



Use **post-interrogation** functions for reporting, root-cause analysis, etc.

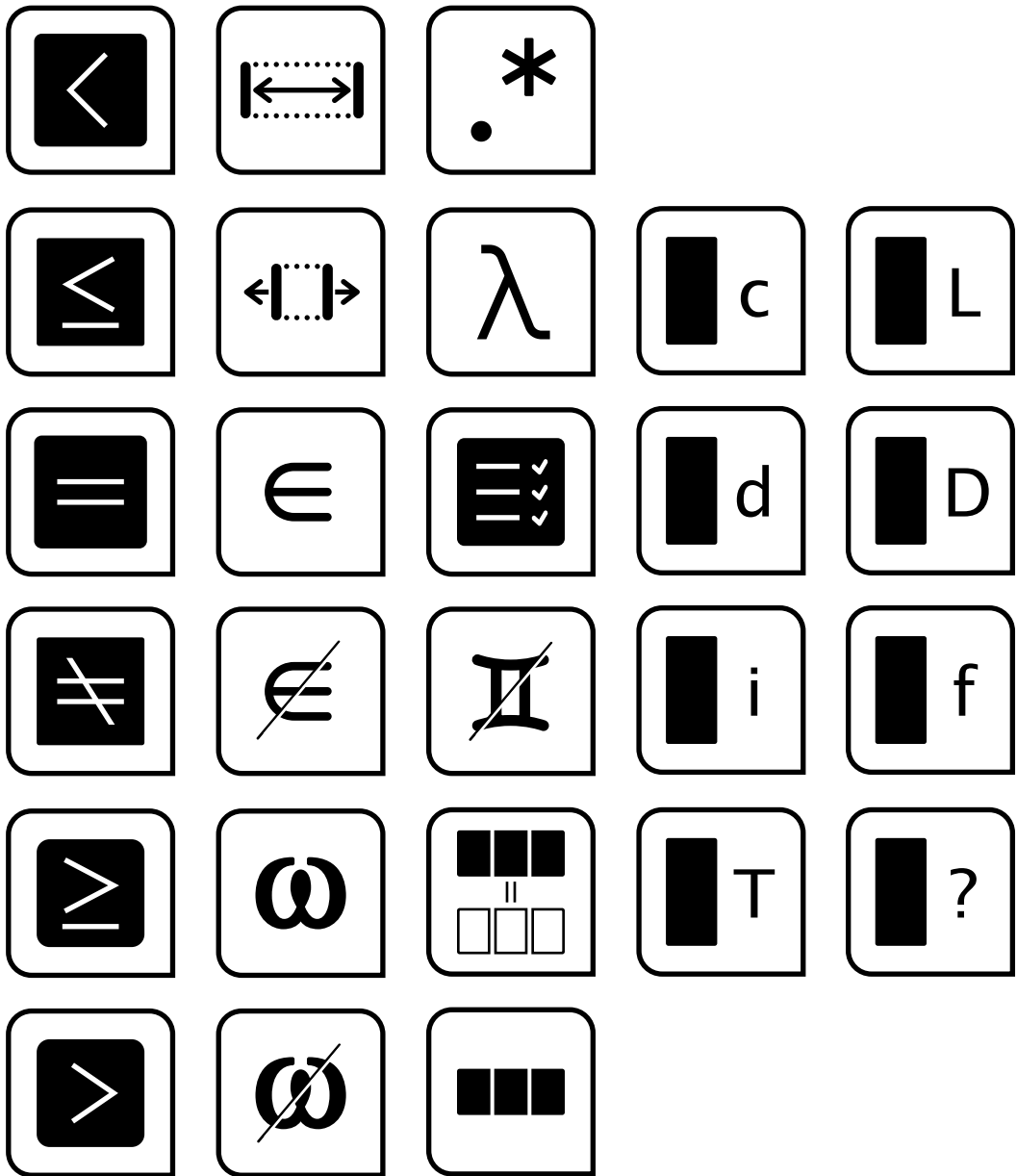
The **pointblank** Data Quality Workflow^①

Understand your data more with a **table scan**.



create_agent()

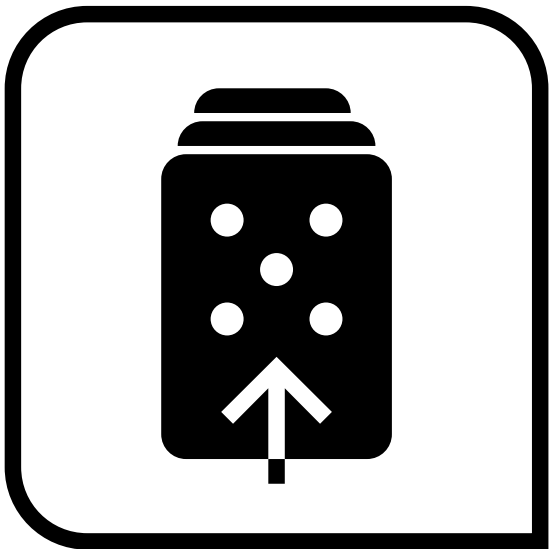
The **agent** is an integral part of the **data quality workflow**.



Send **email** if data quality isn't all that good.

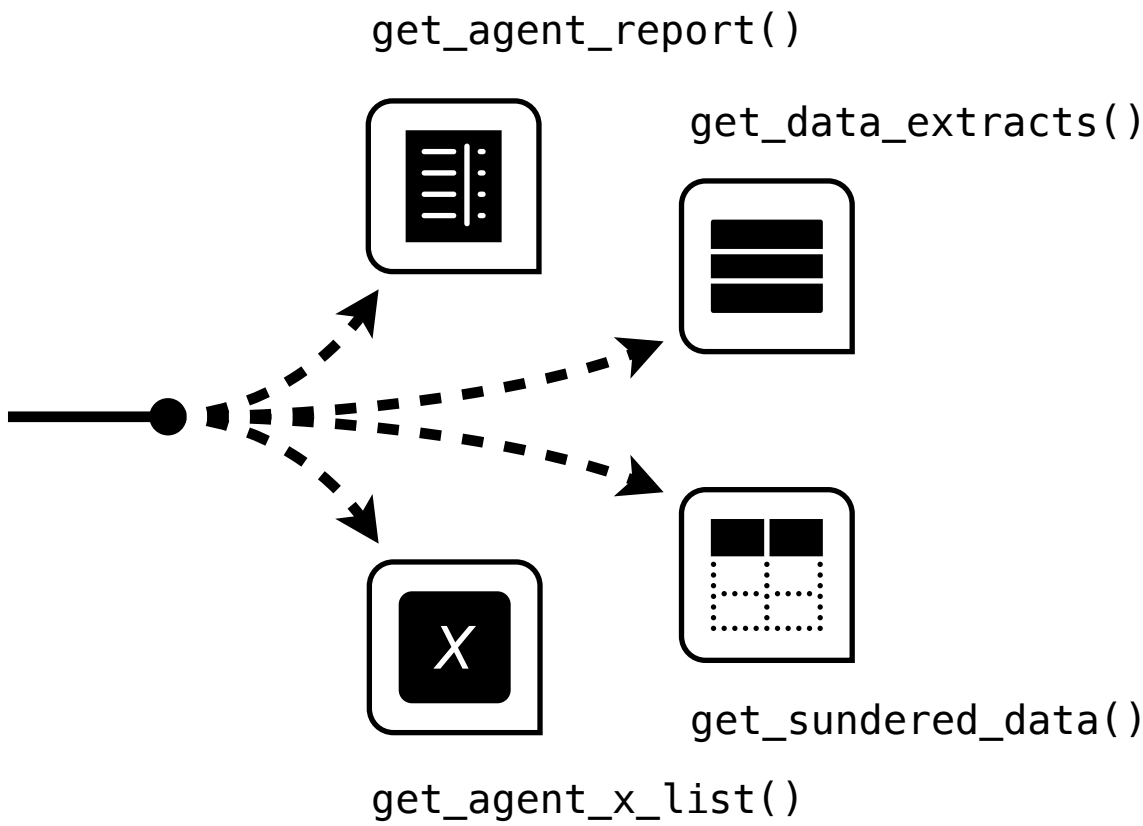


email_blast()



interrogate()

The **agent** performs the **interrogation**. Gathers **intel**.



Use **post-interrogation** functions for reporting, root-cause analysis, etc.

Creating a Validation Plan

a	b	c
yko2	1	23.1
lju7	0	16.3
qib0	1	21.2
sd33	1	24.9
NA	2	NA

Let's start with a simple table

5 rows, 3 columns

Creating a Validation Plan

a	b	c
yko2	1	23.1
lju7	0	16.3
qib0	1	21.2
sd33	1	24.9
NA	2	NA

simple table

5 rows, 3 columns

- 1 All values in **c** should be greater than 15
- 2 All values in **b** should be either 0 or 1
- 3 All values in **a** should fit a pattern of three lowercase letters and a digit
- 4 Values in **c** must be ≥ 20 if **b** is 1; if **b** is 0 then values in **c** must be < 20
- 5 Columns **a**, **b**, and **c** should not have any missing values.

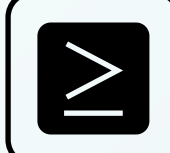
validation plan

5 steps

Creating a Validation Plan

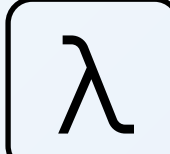
- 1** All values in **c** should be greater than 15
- 2** All values in **b** should be either 0 or 1
- 3** All values in **a** should fit a pattern of three lowercase letters and a digit
- 4** Values in **c** must be ≥ 20 if **b** is 1; if **b** is 0 then values in **c** must be < 20
- 5** Columns **a**, **b**, and **c** should not have any missing values.


validation plan
5 steps

 col_vals_gte()

 col_vals_in_set()

 col_vals_regex()

 col_vals_expr() + case_when()

 col_vals_not_null()

validation functions
5 col_vals_*() functions

Creating a Validation Plan

- 1 All values in **c** should be greater than 15
- 2 All values in **b** should be either 0 or 1
- 3 All values in **a** should fit a pattern of three lowercase letters and a digit
- 4 Values in **c** must be ≥ 20 if **b** is 1; if **b** is 0 then values in **c** must be < 20
- 5 Columns **a**, **b**, and **c** should not have any missing values.

validation plan
5 steps

\geq	<code>col_vals_gte(c, 15)</code>
\in	<code>col_vals_in_set(b, c(0, 1))</code>
\cdot^*	<code>col_vals_regex(a, "[a-z]{3}[0-9]")</code>
λ	<code>col_vals_expr(~ case_when(b == 1 ~ c >= 20, b == 0 ~ c < 20))</code>
\emptyset	<code>col_vals_not_null(vars(a, b, c))</code>

validation functions
5 `col_vals_*`() functions

Interrogating the Table


a	b	c
yko2	1	23.1
lju7	0	16.3
qib0	1	21.2
sd33	1	24.9
NA	2	NA

simple table
5 rows, 3 columns

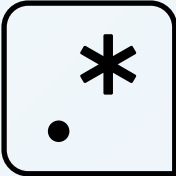
1

 col_vals_gte(c, 15)

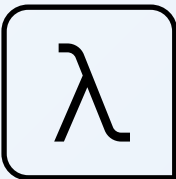
2

 col_vals_in_set(b, c(0, 1))


3

 col_vals_regex(a, "[a-z]{3}[0-9]")

4

 col_vals_expr(~ case_when(
 b == 1 ~ c >= 20,
 b == 0 ~ c < 20))

5

 col_vals_not_null(vars(a, b, c))

validation functions
5 col_vals_*() functions

Interrogating the Table

↓

a	b	c
yko2	1	23.1
lju7	0	16.3
qib0	1	21.2
sd33	1	24.9
NA	2	NA

■

■

■

■

■

INTERROGATION

↓

test units

1

\geq col_vals_gte(c, 15)

REPORT

UNITS	PASS	FAIL
5	4 0.8	1 0.2

Interrogating the Table

↓

a	b		c
yko2	1	<input type="checkbox"/>	23.1
lju7	0	<input type="checkbox"/>	16.3
qib0	1	<input type="checkbox"/>	21.2
sd33	1	<input type="checkbox"/>	24.9
NA	2	<input type="checkbox"/>	NA

test units

2

☐ col_vals_in_set(b, c(0, 1))

Interrogating the Table

↓

a	b		c
yko2	1	■	23.1
lju7	0	■	16.3
qib0	1	■	21.2
sd33	1	■	24.9
NA	2	■	NA

test units

2

\in col_vals_in_set(b, c(0, 1))

REPORT

UNITS	PASS	FAIL
5	4 0.8	1 0.2

Interrogating the Table

↓

a		b	c
yko2	■	1	23.1
lju7	■	0	16.3
qib0	■	1	21.2
sd33	■	1	24.9
NA	■	2	NA

test units

3

`.*` col_vals_regex(a, "[a-z]{3}[0-9]")

REPORT

UNITS	PASS	FAIL
5	3 0.6	2 0.4

Interrogating the Table

a	b	c
yko2	1	23.1
lju7	0	16.3
qib0	1	21.2
sd33	1	24.9
NA	2	NA



test units

4

λ col_vals_expr(~case_when(
b == 1 ~ c >= 20,
b == 0 ~ c < 20))

REPORT

UNITS	PASS	FAIL
4	4	0
	1.0	0

Interrogating the Table

a	b	c
yko2	1	23.1
lju7	0	16.3
qib0	1	21.2
sd33	1	24.9
NA	2	NA

5



col_vals_not_null(vars(a, b, c))

Interrogating the Table

a	b	c
yko2	1	23.1
lju7	0	16.3
qib0	1	21.2
sd33	1	24.9
NA	2	NA

5



col_vals_not_null(vars(a))

6



col_vals_not_null(vars(b))

7



col_vals_not_null(vars(c))

Interrogating the Table

<div>↓</div>		
a		
b		c
yko2	■	123.1
lju7	■	016.3
qib0	■	121.2
sd33	■	124.9
NA	■	2NA

test units

5

☒ col_vals_not_null(vars(a))

6

☒ col_vals_not_null(vars(b))

7

☒ col_vals_not_null(vars(c))

REPORT

UNITS	PASS	FAIL
5	4 0.8	1 0.2

Interrogating the Table

a	b		c
yko2	1	■	23.1
lju7	0	■	16.3
qib0	1	■	21.2
sd33	1	■	24.9
NA	2	■	NA

test units

5



col_vals_not_null(vars(a))

6



col_vals_not_null(vars(b))

7



col_vals_not_null(vars(c))

REPORT

UNITS	PASS	FAIL
5	5	0
	1.0	0

Interrogating the Table

↓

a	b	c
yko2	1	23.1
lju7	0	16.3
qib0	1	21.2
sd33	1	24.9
NA	2	NA

test units

5



col_vals_not_null(vars(a))

6



col_vals_not_null(vars(b))

7



col_vals_not_null(vars(c))

REPORT

UNITS	PASS	FAIL
5	4 0.8	1 0.2

The pointblank Agent Report

	STEP	UNITS	PASS	FAIL
1	col_vals_gte()	5	4 0.8	1 0.2
2	col_vals_in_set()	5	4 0.8	1 0.2
3	col_vals_regex()	5	3 0.6	2 0.4
4	col_vals_expr()	4	4 1.0	0 0
5	col_vals_not_null()	5	4 0.8	1 0.2
6	col_vals_not_null()	5	5 1.0	0 0
7	col_vals_not_null()	5	4 0.8	1 0.2

The pointblank Agent Report

	STEP	UNITS	PASS	FAIL
1	col_vals_gte()	5	4 0.8	1 0.2
2	col_vals_in_set()	5	4 0.8	1 0.2
3	col_vals_regex()	5	3 0.6	2 0.4
4	col_vals_expr()	4	4 1.0	0 0
5	col_vals_not_null()	5	4 0.8	1 0.2
6	col_vals_not_null()	5	5 1.0	0 0
7	col_vals_not_null()	5	4 0.8	1 0.2

For better reporting on data quality, can set thresholds (and even use side effects).

Failure thresholds can be set for three states

W

WARNING

S

STOP

N

NOTIFY

Let's set:

W

to 1

S

to 2

(

N

 not set)

R CODE

1

action_levels(

2

 warn_at = 1,

3

 stop_at = 2

4

)







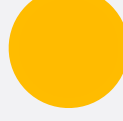




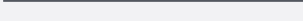


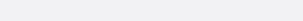


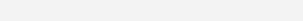


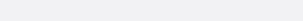
5

6

7

8

The pointblank Agent Report

STEP	UNITS	PASS	FAIL	W	S	N
1 col_vals_gte()	5	4 0.8	1 0.2			
2 col_vals_in_set()	5	4 0.8	1 0.2			
3 col_vals_regex()	5	3 0.6	2 0.4			
4 col_vals_expr()	4	4 1.0	0 0			
5 col_vals_not_null()	5	4 0.8	1 0.2			
6 col_vals_not_null()	5	5 1.0	0 0			
7 col_vals_not_null()	5	4 0.8	1 0.2			

The pointblank Agent Report

Pointblank Validation

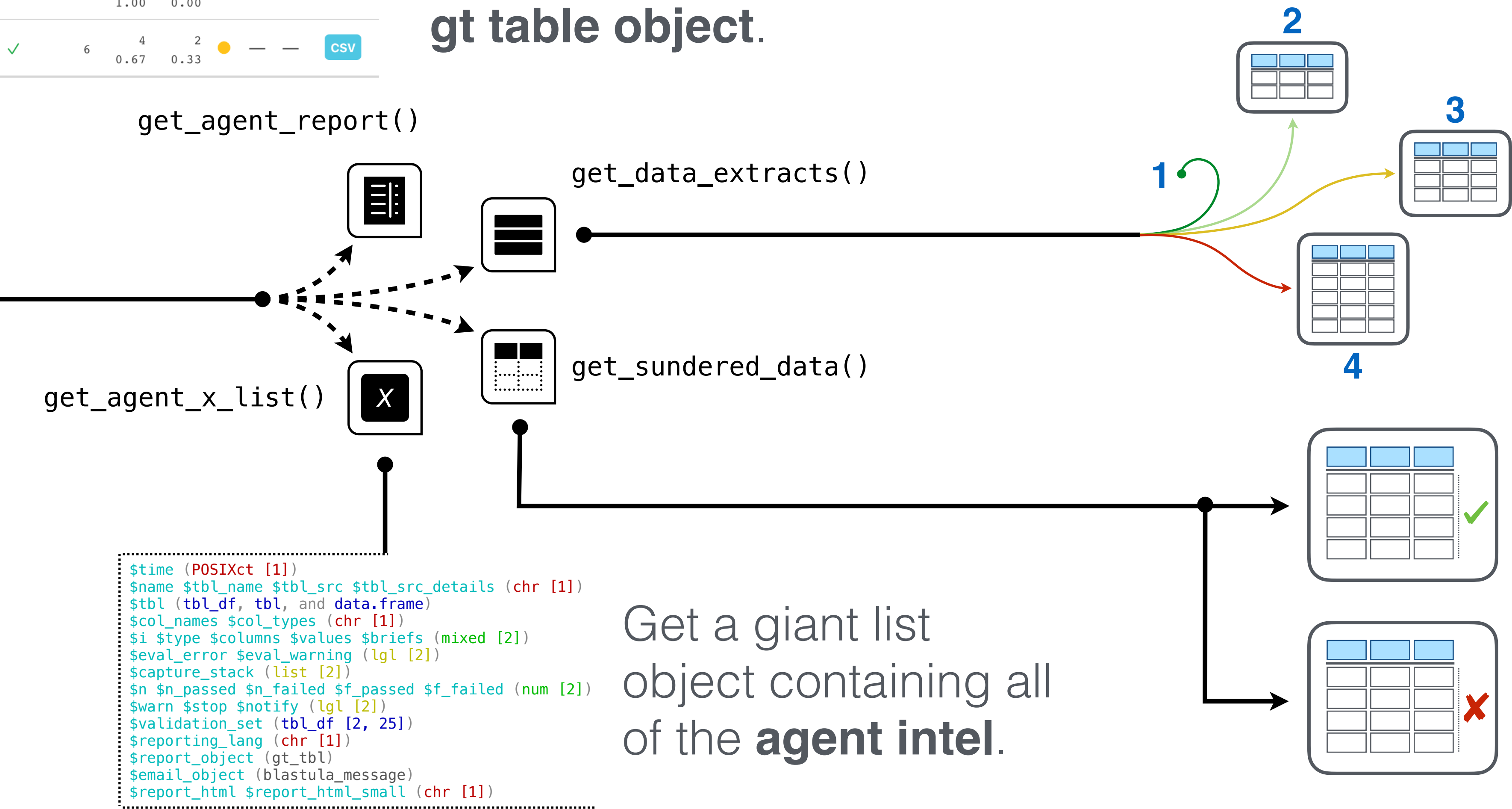
agent_2020-08-13_11:42:20 (2020-08-13 11:42:20)

	STEP	COLUMNS	VALUES	TBL	EVAL	UNITS	PASS	FAIL	W	S	N	EXT	
Yellow	1	col_vals_gte	<div><div></div>c</div>	15	<div>\mathcal{I}</div>	<div>✓</div>	5	<div>40.80</div>	<div>10.20</div>	<div></div>	<div></div>	<div>—</div>	<div>CSV</div>
	2	col_vals_in_set	<div><div></div>b</div>	0, 1	<div>\mathcal{I}</div>	<div>✓</div>	5	<div>40.80</div>	<div>10.20</div>	<div></div>	<div></div>	<div>—</div>	<div>CSV</div>
	3	col_vals_regex	<div><div></div>a</div>	[a-z]{3}[0-9]	<div>\mathcal{I}</div>	<div>✓</div>	5	<div>30.60</div>	<div>20.40</div>	<div></div>	<div></div>	<div>—</div>	<div>CSV</div>
Green	4	col_vals_expr	—	case_when(b ==...	<div>\mathcal{I}</div>	<div>✓</div>	4	<div>41.00</div>	<div>00.00</div>	<div></div>	<div></div>	<div>—</div>	<div>—</div>
Yellow	5	col_vals_not_null	<div><div></div>a</div>	—	<div>\mathcal{I}</div>	<div>✓</div>	5	<div>40.80</div>	<div>10.20</div>	<div></div>	<div></div>	<div>—</div>	<div>CSV</div>
	6	col_vals_not_null	<div><div></div>b</div>	—	<div>\mathcal{I}</div>	<div>✓</div>	5	<div>51.00</div>	<div>00.00</div>	<div></div>	<div></div>	<div>—</div>	<div>—</div>
	7	col_vals_not_null	<div><div></div>c</div>	—	<div>\mathcal{I}</div>	<div>✓</div>	5	<div>40.80</div>	<div>10.20</div>	<div></div>	<div></div>	<div>—</div>	<div>CSV</div>

Other Post-Interrogation Ops

✓	6	6	0	—	—	—	—
✓	6	4	2	—	—	CSV	—

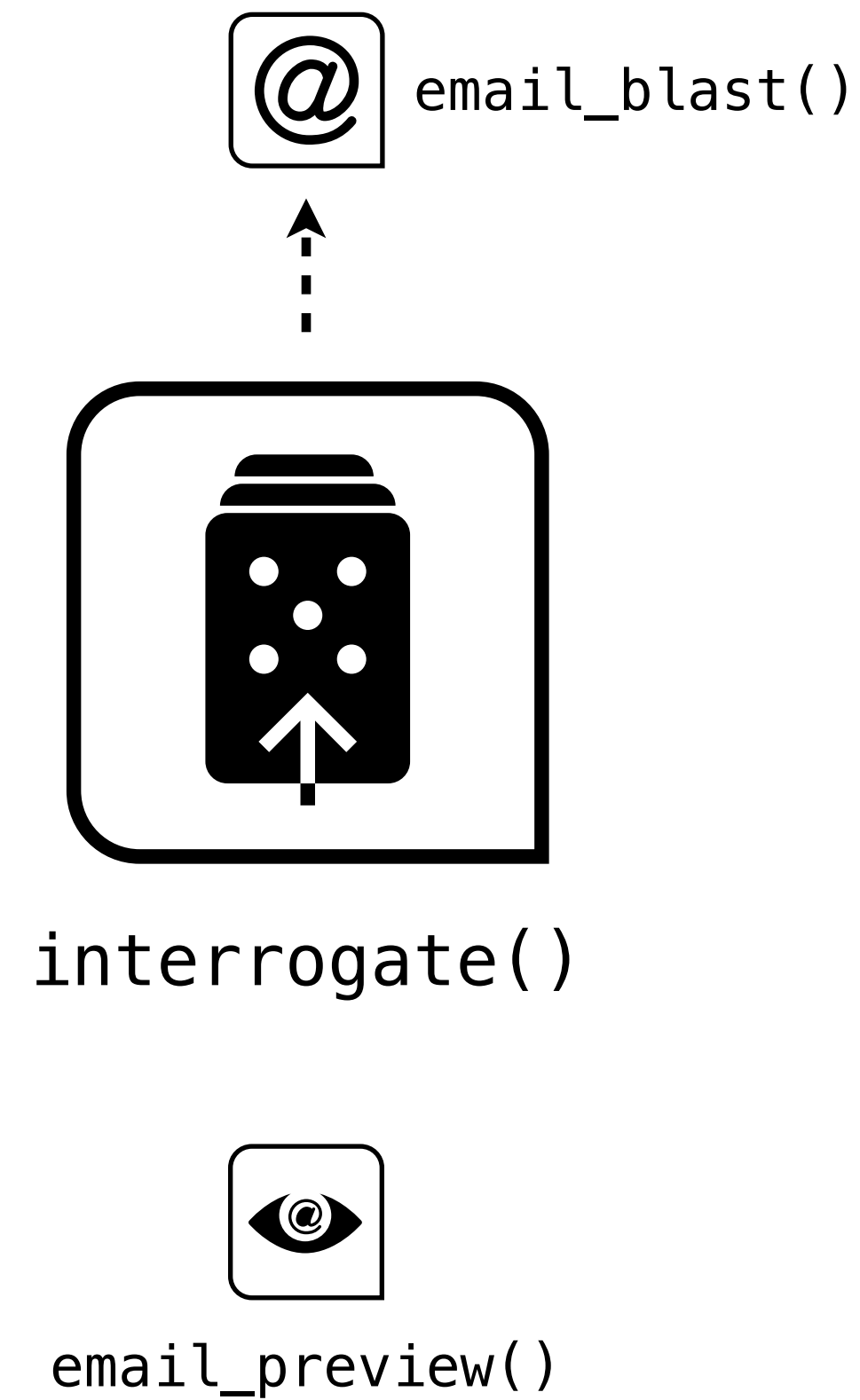
We can get the agent report as a customizable **gt table object**.



```
$time (POSIXct [1])
$name $tbl_name $tbl_src $tbl_src_details (chr [1])
$tbl (tbl_df, tbl, and data.frame)
$col_names $col_types (chr [1])
$i $type $columns $values $briefs (mixed [2])
$eval_error $eval_warning (lgl [2])
$capture_stack (list [2])
$n $n_passed $n_failed $f_passed $f_failed (num [2])
$warn $stop $notify (lgl [2])
$validation_set (tbl_df [2, 25])
$reporting_lang (chr [1])
$report_object (gt_tbl)
$email_object (blastula_message)
$report_html $report_html_small (chr [1])
```


Notifying with Email

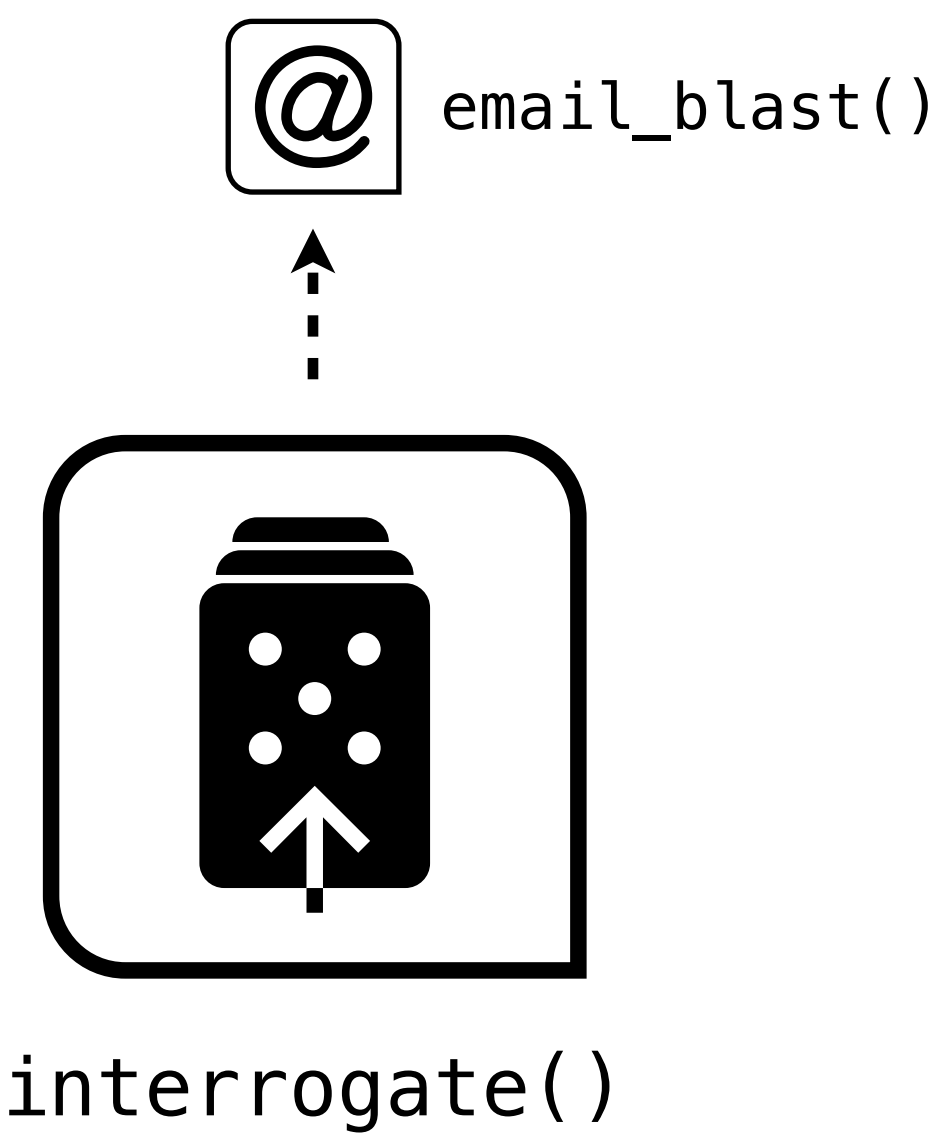
Send email (or *not*) depending on the interrogation results.



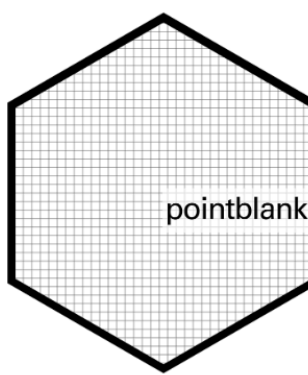
Preview a pointblank email, helpful for customization.

Notifying with Email

Send email (or *not*) depending on the interrogation results.



Preview a pointblank email, helpful for customization.



This **pointblank** validation report, containing 7 validation steps, was initiated on Tuesday, June 23, 2020 at 12:02 PM (EDT).

Pointblank Validation
agent_2020-06-23_12:02:15 (2020-06-23 12:02:16)

	STEP	VALUES	UNITS	PASS	FAIL	W	S	N
1	col_vals_gt	date	13	13 1.00	0 0.00	●	○	—
2	col_vals_gt	g	—	—	—	—	—	—
3	col_vals_regex	[1-9]-[a-z]{3}...	13	13 1.00	0 0.00	●	○	—
4	rows_distinct	—	13	11 0.85	2 0.15	●	○	—
5	col_vals_gt	100	13	13 1.00	0 0.00	●	○	—
6	col_vals_equal	d	13	13 1.00	0 0.00	●	○	—
7	col_vals_between	a, d	13	9 0.69	4 0.31	●	●	—

©

Validation performed via the pointblank R package.

INFORMATION AND PACKAGE DOCUMENTATION

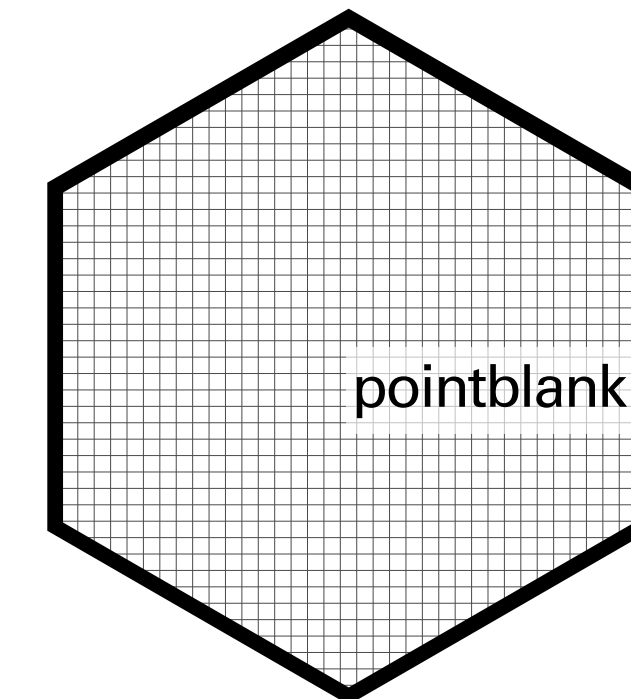
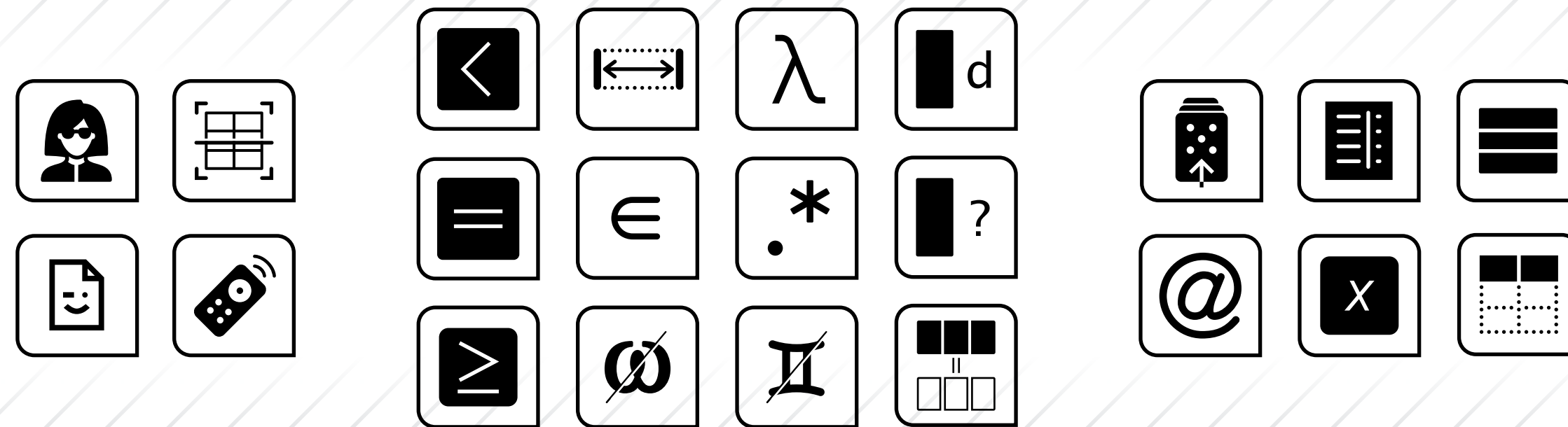
This is the stock email w/o any customization.

It's made with the **blastula** package. It's an HTML email that is well-tested in multiple email clients.



Demo

The **pointblank** R Package



github.com/rich-iannone/pointblank

github.com/rich-iannone/presentations



rich-iannone



@riannone



rich@rstudio.com