

Package: duckspending (via r-universe)

May 31, 2026

Title Access to USAspending Data in a DuckLake

Version 0.1.0

Description Provides a connection to USAspending data in a DuckLake, allowing users to perform bulk analyses difficult to achieve with the USAspending API.

License MIT + file LICENSE

Encoding UTF-8

Roxygen list(markdown = TRUE)

RoxygenNote 8.0.0

Imports cli, DBI, duckdb, lifecycle, methods, rlang, tibble, utils, vctrs

Suggests connections, dbplyr, dplyr, jsonlite, knitr, labelled, rmarkdown, rscontract, testthat (>= 3.0.0), withr, pkgdown, rprojroot

VignetteBuilder knitr

Config/testthat/edition 3

URL <https://usaspending.grant-witness.us/r-package/>

Contact info@grant-witness.us

Config/pak/sysreqs xz-utils

Repository <https://noamross.r-universe.dev>

Date/Publication 2026-05-20 00:13:56 UTC

RemoteUrl <https://usaspending.grant-witness.us/r-package/duckspending>

RemoteRef HEAD

RemoteSha fd2dfb428242dad1aa9cd63c77217afcfa4a9269

Contents

| | |
|---|---|
| connection_view.duckspending_connection | 2 |
| duckspending_attach | 3 |

| | |
|--------------------------------------|----|
| duckspending_attached | 3 |
| duckspending_clear_cache | 4 |
| duckspending_connection | 4 |
| duckspending_describe | 6 |
| duckspending_descriptions | 7 |
| duckspending_detach | 8 |
| duckspending_label | 8 |
| duckspending_label_default | 9 |
| duckspending_latest | 10 |
| duckspending_snapshots | 10 |
| duckspending_table_helpers | 11 |
| duckspending_tables | 13 |
| duckspending_tbl | 14 |
| duckspending_view | 15 |

Index 16

connection_view.duckspending_connection
IDE Connections pane integration for a duckspending connection

Description

S3 method for `connections::connection_view()`. Builds an `rscontract::rscontract_spec()` that wires the pane's drill-down to `pane_object_list()` and `pane_object_columns()` so the IDE sees a filtered, user-facing view of the DuckLake.

Usage

```
connection_view.duckspending_connection(
  con,
  connection_code = "",
  host = "",
  name = "",
  connection_id = ""
)
```

Arguments

`con` A `duckspending_connection`.
`connection_code` Optional connection code shown by the IDE.
`host, name, connection_id` See `connections::connection_view()`.

Value

The result of `rscontract::rscontract_open()`, invisibly.

`duckspending_attach` *Attach an additional snapshot to an open connection*

Description

Caches the snapshot's catalog locally (if needed), ATTACHes it under a hidden alias, and creates a top-level schema named after the date with one view per user table. After this call the snapshot is queryable as `<conn>.<YYYYMMDD>.<table>`.

Usage

```
duckspending_attach(snapshot, conn = duckspending_connection())
```

Arguments

`snapshot` Snapshot identifier (YYYYMMDD character, `Date`, or numeric).
`conn` A `duckspending_connection`. Defaults to the cached default `duckspending_connection()`.

Details

If `conn` is the package's cached default connection (the one `duckspending_connection()` returns with no args), the cache entry is updated in place so subsequent default-arg callers see the new state.

Value

`conn`, invisibly, with the new snapshot in its metadata.

`duckspending_attached`
Currently-attached snapshots on a connection

Description

Currently-attached snapshots on a connection

Usage

```
duckspending_attached(conn = duckspending_connection())
```

Arguments

`conn` A `duckspending_connection`. Defaults to the cached default `duckspending_connection()`.

Value

Character vector of YYYYMMDD dates, newest first.

`duckspending_clear_cache`

Clear the on-disk HTTP cache for a duckspending connection

Description

Calls `cache_httpfs_clear_cache()` on the underlying DuckDB session to delete all cached parquet blocks from disk. Useful when you want to force fresh fetches or reclaim disk space.

Usage

```
duckspending_clear_cache(conn = duckspending_connection())
```

Arguments

`conn` A `duckspending_connection`. Defaults to the cached default `duckspending_connection()`.

Details

Has no effect (with a message) if the connection was opened with `cache = FALSE` or if the `cache_httpfs` extension failed to load.

Value

`conn`, invisibly.

`duckspending_connection`

Connect to the USAspending DuckLake archive

Description

Opens a DuckDB connection with one or more monthly DuckLake snapshots attached. Each attached snapshot is exposed as a top-level schema named after its date ("20260206", "20260306", ...). A `latest` schema mirrors the newest attached snapshot for convenience.

Usage

```
duckspending_connection(
  snapshots = "latest",
  base_url = default_base_url(),
  catalog_dir = tools::R_user_dir("duckspending", which = "cache"),
  download_catalog = TRUE,
  open_pane = TRUE,
  cache = .duckspending_cache_default(),
  new_conn = FALSE
)
```

Arguments

| | |
|-------------------------------|--|
| <code>snapshots</code> | Snapshot selector. One of: <ul style="list-style-type: none"> • "latest" (default): attach only the newest snapshot per <code>duckspending_snapshots()</code>. • "all": attach every snapshot in the manifest. • A character vector of YYYYMMDD strings, or coercible (<code>Date</code> or numeric), naming the snapshots to attach. |
| <code>base_url</code> | Base URL of the public archive. Override for testing via <code>options(duckspending.base_url=)</code> or the <code>DUCKSPENDING_BASE_URL</code> env var; defaults to https://usaspending.grant-witness.us . |
| <code>catalog_dir</code> | Local cache directory for snapshot catalog files. |
| <code>download_catalog</code> | Logical; when <code>TRUE</code> (default) catalog files are cached locally before attaching, when <code>FALSE</code> the catalog is attached directly over HTTP. |
| <code>open_pane</code> | Logical; when <code>TRUE</code> (default) and the <code>connections</code> package is available, register the connection with the IDE Connections pane. |
| <code>cache</code> | Logical; when <code>TRUE</code> (default) the <code>cache_httpfs</code> community extension is loaded to persist fetched parquet blocks on disk, speeding up repeated and parallel queries. Override via <code>options(duckspending.cache = FALSE)</code> or <code>DUCKSPENDING_CACHE=false</code> . Set to <code>FALSE</code> to skip caching. |
| <code>new_conn</code> | Logical; when <code>TRUE</code> force a new connection even if a cached default exists. |

Details

By default a single snapshot — the newest one announced by the public manifest — is attached. To attach a specific set, pass `snapshots = ...`.

The first call with default arguments caches its connection in a package-level environment; subsequent default calls reuse it, so all the other `duckspending_*` helpers that default `conn = duckspending_connection()` share a single underlying DuckDB session. Calls with non-default arguments are not cached.

Value

A `duckspending_connection` (`duckdb_connection` subclass).

Examples

```
## Not run:
con <- duckspending_connection()
duckspending_tbl("subaward_search") # uses the cached default
DBI::dbDisconnect(con, shutdown = TRUE)

## End(Not run)
```

 duckspending_describe

Describe snapshots, tables, and columns visible on a connection

Description

Three modes:

- With no `table`, returns one row per user-facing table in each attached snapshot (plus the `latest` alias).
- With a `table`, returns a column-level summary for that table in the requested `snapshot` (defaults to `latest`).

Usage

```
duckspending_describe(
  table = NULL,
  snapshot = NULL,
  descriptions = duckspending_descriptions(),
  include_source = FALSE,
  conn = duckspending_connection()
)
```

Arguments

| | |
|-----------------------------|--|
| <code>table</code> | Optional table name. When <code>NULL</code> , lists snapshot/table pairs. |
| <code>snapshot</code> | Snapshot selector for the column-level mode. <code>NULL</code> uses the <code>latest</code> schema. |
| <code>descriptions</code> | Optional list returned by <code>duckspending_descriptions()</code> . When supplied, populates the <code>description</code> column. |
| <code>include_source</code> | If <code>TRUE</code> , the column-level mode also returns a <code>source</code> column attributing each description (e.g. "Data Dictionary Crosswalk", "Inferred"). Defaults to <code>FALSE</code> to keep the common case tidy. |
| <code>conn</code> | A <code>duckspending_connection</code> . Defaults to the cached default <code>duckspending_connection()</code> . |

Details

Internal bookkeeping (DuckLake metadata, hidden `_dl_*` databases, `_snapshots` marker, DuckDB metadata schemas) is filtered out.

Value

A tibble.

```
duckspending_descriptions
```

Fetch and parse a descriptions JSON file for a USAspending DuckLake

Description

Descriptions live in a JSON file at the bucket root next to the catalog (default URL is derived from the production catalog URL). When the file is unreachable (e.g. not yet published, transient network error), returns NULL so callers can fall back to schema-only output.

Usage

```
duckspending_descriptions(
  url = NULL,
  dir = tools::R_user_dir("duckspending", which = "cache"),
  refresh = FALSE
)
```

Arguments

| | |
|----------------------|--|
| <code>url</code> | URL or local path to a descriptions JSON. Defaults to <code><base_url>/descriptions.json</code> , where <code>base_url</code> is the active archive URL (<code>default_base_url()</code>). |
| <code>dir</code> | Directory to cache the file in. |
| <code>refresh</code> | Force a re-check this session (ignores the per-session memoisation). The conditional ETag/Last-Modified check still applies, so a still-fresh remote file isn't re-downloaded. |

Details

The file is cached under `tools::R_user_dir()` and refreshed conditionally via the ETag / Last-Modified headers: the first call per R session sends a HEAD request, and a GET only follows when the remote validators differ from the cached ones. Subsequent calls in the same session reuse the cached parse without any HTTP traffic.

If the per-session check fails (e.g. no network), the failure is silent: a cached copy is returned when available, NULL otherwise, and the next call within the same session re-attempts the check — so a transient outage doesn't permanently disable the help-text refresh.

Expected JSON shape:

```
{
  "tables": {
    "<schema>.<table>": {
      "description": "...",
      "columns": {"<column>": "..."}
    }
  }
}
```

Value

A parsed list, or NULL if the file could not be fetched.

`duckspending_detach` *Detach a snapshot from an open connection*

Description

Drops the user-facing schema ("YYYYMMDD") and DETACHes the hidden underlying catalog. Removes the snapshot from the connection metadata. Refuses to detach the snapshot currently bound to `latest`.

Usage

```
duckspending_detach(snapshot, conn = duckspending_connection())
```

Arguments

`snapshot` Snapshot identifier.

`conn` A `duckspending_connection`. Defaults to the cached default `duckspending_connection()`.

Value

`conn`, invisibly.

`duckspending_label` *Apply column descriptions as labels to a data frame*

Description

Stamps `attr(col, "label")` on every column in `data` whose name has a matching description in the USAspending descriptions dictionary. Both Positron's Data Explorer and RStudio's `View()` honour `attr(., "label")` and surface it under the column header — the same convention used by the `labelled` and `haven` packages.

Usage

```
duckspending_label(
  data,
  table = NULL,
  descriptions = duckspending_descriptions()
)
```

Arguments

| | |
|---------------------------|--|
| <code>data</code> | A data frame or tibble. |
| <code>table</code> | Optional table name. When supplied, only descriptions for that table are used (unambiguous; recommended when you know the source). When <code>NULL</code> , descriptions from all tables are merged into a flat lookup; on collisions (same column name in multiple tables), the first match wins. |
| <code>descriptions</code> | Optional list returned by <code>duckspending_descriptions()</code> . <code>NULL</code> makes this a no-op (useful when the dictionary is unreachable). |

Details

Columns whose names aren't in the dictionary are left untouched.

Value

`data`, with `attr(, "label")` set on matching columns.

`duckspending_label_default`

Resolve whether labels should be applied by default

Description

Checks (in order): `options(duckspending.label = ...)`, the `DUCKSPENDING_LABEL` env var, and falls back to `TRUE`. The env var is treated as falsy if set to `"false"`, `"0"`, `"no"`, `"off"` (case-insensitive); anything else (including unset-but-checked) is truthy.

Usage

```
duckspending_label_default()
```

Value

A single logical.

`duckspending_latest` *The newest available snapshot date.*

Description

Convenience equivalent of `duckspending_snapshots()[1]`. Returns `NA_character_` if no manifest is reachable.

Usage

```
duckspending_latest(
  refresh = FALSE,
  base_url = default_base_url(),
  cache_dir = tools::R_user_dir("duckspending", which = "cache")
)
```

Arguments

| | |
|------------------------|--|
| <code>refresh</code> | Logical; bypass the local cache and re-fetch. |
| <code>base_url</code> | Base URL (default: production). Override with <code>options(duckspending.base_url = ...)</code> for testing. |
| <code>cache_dir</code> | Local cache directory. |

Value

YYYYMMDD character scalar or `NA_character_`.

`duckspending_snapshots`
List available USAspending snapshots

Description

Fetches the small text manifest at `{base_url}/snapshots.txt` and returns the date strings, newest first. The manifest is cached in `tools::R_user_dir("duckspending", "cache")` for one hour to keep repeat calls cheap.

Usage

```
duckspending_snapshots(
  refresh = FALSE,
  base_url = default_base_url(),
  cache_dir = tools::R_user_dir("duckspending", which = "cache"),
  ttl_seconds = 3600L
)
```

Arguments

| | |
|--------------------------|--|
| <code>refresh</code> | Logical; bypass the local cache and re-fetch. |
| <code>base_url</code> | Base URL (default: production). Override with <code>options(duckspending.base_url = ...)</code> for testing. |
| <code>cache_dir</code> | Local cache directory. |
| <code>ttl_seconds</code> | Cache TTL in seconds (default 1 hour). |

Value

A character vector of YYYYMMDD dates, newest first. Empty character if the manifest is not found.

duckspending_table_helpers

Convenience accessors for the most-used USAspending tables

Description

Each function below returns a non-collected `dplyr` table reference (i.e. a `tbl_lazy`) for one of the headline tables in the archive. They are thin wrappers around `duckspending_tbl()` that exist so the table names users care about most are exported, autocompletable, and discoverable via `?duckspending` — you don't have to remember the exact column-name spelling of `subaward_search` etc.

Usage

```

duckspending_subawards(
  snapshot = NULL,
  label = duckspending_label_default(),
  conn = duckspending_connection()
)

duckspending_assistance_transactions(
  snapshot = NULL,
  label = duckspending_label_default(),
  conn = duckspending_connection()
)

duckspending_procurement_transactions(
  snapshot = NULL,
  label = duckspending_label_default(),
  conn = duckspending_connection()
)

duckspending_agencies(
  snapshot = NULL,

```

```

    label = duckspending_label_default(),
    conn = duckspending_connection()
  )

  duckspending_recipients(
    snapshot = NULL,
    label = duckspending_label_default(),
    conn = duckspending_connection()
  )

```

Arguments

| | |
|-----------------------|---|
| <code>snapshot</code> | Optional snapshot selector (YYYYMMDD, Date, or numeric); defaults to the connection's <code>latest</code> schema. |
| <code>label</code> | Whether to attach column descriptions as <code>attr(, "label")</code> when the result is <code>collect()</code> ed. Defaults to <code>duckspending_label_default()</code> . |
| <code>conn</code> | A <code>duckspending_connection</code> . Defaults to the cached default <code>duckspending_connection()</code> . |

Details

Everything is lazy: build a `dplyr` pipeline and call `dplyr::collect()` when you want the result in R, or run `dplyr::show_query()` to inspect the SQL. To use the legacy raw tables not listed here, call `duckspending_tbl()` directly.

Value

A `tbl_lazy` reference.

Functions

- `duckspending_subawards()`: Sub-award table (`subaward_search`) — one row per disclosed sub-award.
- `duckspending_assistance_transactions()`: Federal Assistance (FABS) transactions — `source_assistance_transaction`.
- `duckspending_procurement_transactions()`: Federal Procurement (FPDS) transactions — `source_procurement_transaction`.
- `duckspending_agencies()`: Toptier agencies (`toptier_agency`) — one row per department/agency at the top of the federal hierarchy.
- `duckspending_recipients()`: Recipient profile table (`recipient_profile`) — one row per known funding recipient with rolled-up totals.

duckspending_tables *Tables published in the USAspending DuckLake archive*

Description

This help topic is rendered dynamically from `<base>/descriptions.json` on the public archive, so it always reflects the descriptions currently shipping alongside the catalog – no package release is required when descriptions change.

Details

The table below lists every user-facing table in the archive together with a short description. Column-level descriptions are available via `duckspending_describe()`.

If the descriptions JSON is unreachable (network outage, file not yet published) the table will say so and you can still inspect the schema interactively with `duckspending_describe()`.

Table

agency
 appropriation_account_balances
 award_category
 budget_authority
 bureau_title_lookup
 c_to_d_linkage_updates
 cgac
 dabs_submission_window_schedule
 disaster_emergency_fund_code
 duns
 external_data_load_date
 external_data_type
 federal_account
 financial_accounts_by_awards
 financial_accounts_by_program_activity_object_class
 frec
 frec_map
 gtas_sf133_balances
 historic_parent_duns
 historical_appropriation_account_balances
 naics
 object_class
 office
 overall_totals
 parent_award
 psc
 recipient_profile
 ref_city_county_state_code
 ref_country_code

Description

Joining table that pairs each subtier agency with its top-level agency.
 DATA Act File A: TAS-level totals for each agency reporting period.
 Reference table mapping internal award type/category codes to human-readable descriptions.
 Aggregated agency-level budget authority figures by fiscal year.
 Crosswalk of agency/bureau codes to human-readable bureau names.
 Manual corrections to the linkage between File C financial accounts and agency codes.
 Common Government-wide Accounting Classification (CGAC) codes.
 Schedule of DATA Act reporting periods: open/close dates.
 Disaster Emergency Fund Codes (DEFC) — tag obligations.
 Crosswalk between legacy DUNS numbers and current USAspending DUNS numbers.
 Tracks when each external feed (SAM.gov, FPDS, FABS) was last updated.
 Reference table naming the external data sources tracked in the archive.
 Federal Account groupings — the parent level above TAAs.
 DATA Act File C: the bridge between budget appropriations and agency codes.
 DATA Act File B: how each TAS spent money, broken down by object class.
 Federal Reporting Entity Codes — used by certain Treasury entities.
 FREC-to-CGAC reconciliation table used to map FREC codes to CGAC codes.
 GTAS/SF-133 government-wide budget execution balances.
 Legacy parent-DUNS crosswalk for historical recipient matching.
 Pre-DATA Act TAS-level account balances, sourced from the USAspending archive.
 North American Industry Classification System codes — used for NAICS reporting.
 Object class codes (OMB Circular A-11) — classify federal spending.
 Contracting offices and sub-organizations within subtier agencies.
 Top-line fiscal-year aggregates (total budget authority) used for reporting.
 IDV (Indefinite Delivery Vehicle) and ordering-agreement codes.
 Product and Service Codes — more granular than NAICS codes.
 Recipient-level summary records used for recipient search and reporting.
 City × county × state FIPS code reference used for placename reporting.
 ISO and USAspending country codes used in place-of-payment reporting.

| | |
|--------------------------------|--|
| ref_population_cong_district | ACS-derived population estimates by state × congressional district. |
| ref_population_county | ACS-derived population estimates by state × county. Updated quarterly. |
| ref_program_activity | Program activity codes — identify the specific activity of a transaction. |
| references_cfda | Federal Assistance Listings (formerly CFDA — Catalog of Federal Domestic Assistance). |
| references_definition | Glossary table with plain-language and official definitions of program activity codes. |
| source_assistance_transaction | Assistance transaction actions (grants, loans, direct payments). |
| source_procurement_transaction | Contract transaction actions (FPDS) — one row per contract. |
| state_data | State-level metadata: FIPS codes, abbreviations, and names. |
| subaward_search | Subaward records reported by prime recipients to FRS (Federal Reporting System). |
| submission_attributes | One row per DATA Act submission. The authoritative source of program activity codes. |
| subtier_agency | Subtier agencies — bureaus, administrations, and offices of executive agencies. |
| toptier_agency | Cabinet-level agencies (Department of Defense, Department of Justice, etc.). |
| treasury_appropriation_account | Treasury Account Symbol (TAS) definitions — the atomic source of program activity codes. |
| zips_grouped | ZIP-to-county/congressional-district crosswalk used by geographers. |

See Also

[duckspending_describe\(\)](#) for live, column-level details on a connected DuckLake; [duckspending_descriptions\(\)](#) for the raw parsed JSON.

`duckspending_tbl` *Lazy table accessor for a USAspending DuckLake snapshot*

Description

Returns a non-collected `dplyr` table reference to `<snapshot>.<name>`, where `<snapshot>` defaults to the connection's `latest` schema. Verifies the table exists first so users get a clear error instead of a deferred SQL failure.

Usage

```
duckspending_tbl(
  name,
  snapshot = NULL,
  label = duckspending_label_default(),
  conn = duckspending_connection()
)
```

Arguments

| | |
|-----------------------|---|
| <code>name</code> | Table name (no quoting needed). |
| <code>snapshot</code> | Optional snapshot selector. One of: <ul style="list-style-type: none"> • <code>NULL</code> (default): use the connection's <code>latest</code> schema. • <code>"latest"</code>: explicit equivalent of <code>NULL</code>. • <code>YYYYMMDD</code> character / <code>Date</code> / numeric naming an attached snapshot. |

| | |
|-------|---|
| label | Whether to attach column descriptions as <code>attr(., "label")</code> when the lazy table is <code>collect()</code> ed. Defaults to <code>duckspending_label_default()</code> (controlled by <code>options(duckspending.label)</code> and <code>DUCKSPENDING_LABEL</code>). |
| conn | A <code>duckspending_connection</code> . Defaults to the cached default <code>duckspending_connection()</code> . |

Value

A `tbl_lazy` (dbplyr) reference to the table.

`duckspending_view` *Open the IDE Connections pane for a duckspending connection*

Description

Thin convenience wrapper around `connections::connection_view()` that dispatches to the `duckspending_connection` S3 method. Useful after creating a connection with `open_pane = FALSE`, or to re-open the pane after closing it.

Usage

```
duckspending_view(conn = duckspending_connection())
```

Arguments

conn A `duckspending_connection`. Defaults to the cached default `duckspending_connection()`.

Value

conn, invisibly.

Index

- * **documentation**
 - duckspending_tables, 13

- connection_view.duckspending_connection, 2
- connections::connection_view(), 2, 15

- default_base_url(), 7
- dplyr::collect(), 12
- dplyr::show_query(), 12
- duckspending_agencies
 - (*duckspending_table_helpers*), 11
- duckspending_assistance_transactions
 - (*duckspending_table_helpers*), 11
- duckspending_attach, 3
- duckspending_attached, 3
- duckspending_clear_cache, 4
- duckspending_connection, 4
- duckspending_connection(), 3, 4, 6, 8, 12, 15
- duckspending_describe, 6
- duckspending_describe(), 13, 14
- duckspending_descriptions, 7
- duckspending_descriptions(), 6, 9, 14
- duckspending_detach, 8
- duckspending_label, 8
- duckspending_label_default, 9
- duckspending_label_default(), 12, 15
- duckspending_latest, 10
- duckspending_procurement_transactions
 - (*duckspending_table_helpers*), 11
- duckspending_recipients
 - (*duckspending_table_helpers*), 11
- duckspending_snapshots, 10
- duckspending_snapshots(), 5

- duckspending_subawards
 - (*duckspending_table_helpers*), 11
- duckspending_table_helpers, 11
- duckspending_tables, 13
- duckspending_tbl, 14
- duckspending_tbl(), 11, 12
- duckspending_view, 15

- pane_object_columns(), 2
- pane_object_list(), 2

- rscontract::rscontract_open(), 2
- rscontract::rscontract_spec(), 2

- tools::R_user_dir(), 7