# Batch Time Analysis of Transactional Data

## Description

Lenodo is a multinational e-commerce organization that sells products directly to consumers. The database administrator exports the data every night in a CSV file, but this export functionality is unused. Lenodo wants to use this data to uncover insights about the most-sold item and the countries where customers have bought this item.

You are a data analytics consultant, and you're asked to provide valuable insights and statistics across products, brands, categories, segments to the marketing, product, sales, and procurement teams and inform them about which product has the highest amount of sales and which product and its marketing needs the most improvement. These statistics will help to run effective digital marketing campaigns. The scope of this project is limited to data engineering and analysis.

**Objective:**

To use AWS Big Data stack for data engineering to analyze transactions, uncover patterns, and share actionable insights

**Steps to perform:**

1. Create an S3 bucket with a unique name and upload the CSV file to the S3 bucket (ensure that the file is in UTF-8 format only)
2. Create a crawler to crawl the CSV data and generate a metadata catalog
3. Create a Glue job to transform the data into the Parquet format as CSV is not optimal for data warehouse queries
4. Add another crawler to crawl the Parquet data files to generate the metadata catalog of the Parquet file in order to query it with Athena
5. Query the data to identify the best-selling item and countries where customers have bought the most-sold item using Athena

## 1. Setup AWS S3 Bucket:

Objective: Store the CSV data file securely and reliably in the cloud.

Action Items:

Go to the AWS Management Console, navigate to S3, and create a new bucket.

Ensure the name is unique, follows AWS naming conventions, and is region-appropriate for your analysis needs.

Upload the CSV file to the newly created S3 bucket. Make sure the CSV file is encoded in UTF-8 to avoid any compatibility issues.

### Create bucket Info

Buckets are containers for data stored in S3.

**General configuration**

AWS Region

US East (N. Virginia) us-east-1 ▼

Bucket type   Info

● General purpose
Recommended for most use cases and access patterns. General purpose buckets are the original S3 bucket type. They allow a mix of storage classes that redundantly store objects across multiple Availability Zones.

○ Directory - *New*
Recommended for low-latency use cases. These buckets use only the S3 Express One Zone storage class, which provides faster processing of data within a single Availability Zone.

Bucket name   Info

project1bysamir

Bucket name must be unique within the global namespace and follow the bucket naming rules. See rules for bucket naming 🔗

Copy settings from existing bucket - *optional*
Only the bucket settings in the following configuration are copied.

**Choose bucket**

Format: s3://bucket/prefix

---

Amazon S3 > Buckets > project1bysamir

### project1bysamir Info

Objects | Properties | Permissions | Metrics | Management | Access Points

**Objects (1)** Info    [C]    [Copy S3 URI]    [Copy URL]    [Download]    [Open 🔗]    [Delete]    [Actions ▼]    [Create folder]    [Upload]
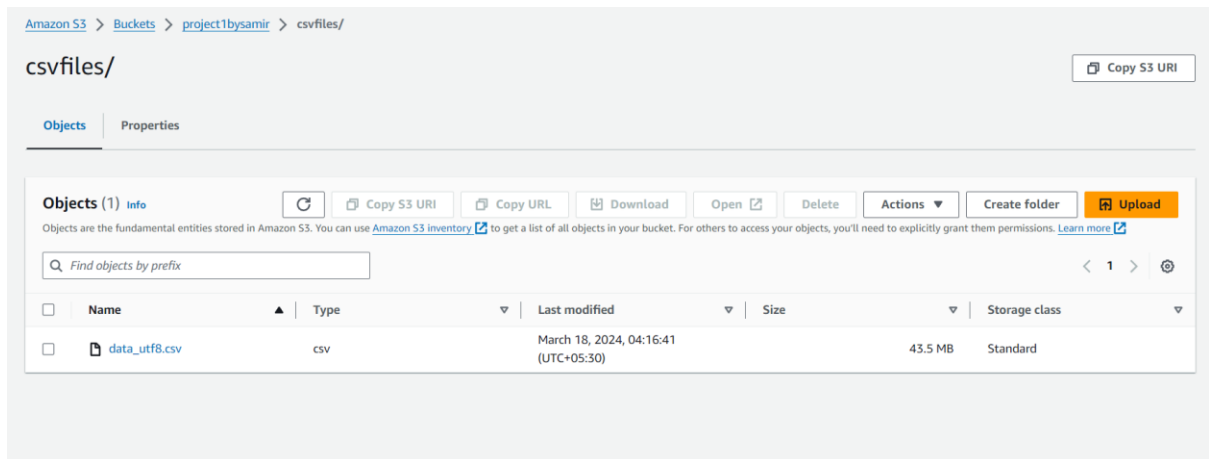
Objects are the fundamental entities stored in Amazon S3. You can use Amazon S3 inventory 🔗 to get a list of all objects in your bucket. For others to access your objects, you'll need to explicitly grant them permissions. Learn more 🔗

🔍 Find objects by prefix                                                                                      < 1 >   ⚙

| ☐ | Name | ▲ | Type | ▽ | Last modified | ▽ | Size | ▽ | Storage class | ▽ |
|---|---|---|---|---|---|---|---|---|---|---|
| ☐ | 📁 csvfiles/ | | Folder | | - | | - | | - | |

## 2. Create and Run AWS Glue Crawler for CSV Data:

Objective: Automatically discover and catalog metadata from the CSV data in S3.

Action Items:

In AWS Glue, create a new crawler. Set the data store to point to the S3 bucket where the CSV file is stored.

Choose an IAM role that has permissions to access both AWS Glue services and the S3 bucket.

Configure the crawler to run on demand or schedule it as needed. Once configured, run the crawler.

Upon completion, the crawler will create a database and table(s) in the AWS Glue Data Catalog, representing the structure of your CSV data.

# Choose data sources and classifiers

## Data source configuration

Is your data already mapped to Glue tables?

○ **Not yet**
Select one or more data sources to be crawled.

○ Yes
Select existing tables from your Glue Data Catalog.

**Data sources (1)** Info

[Edit] [Remove] [Add a data source]

The list of data sources to be scanned by the crawler.

| | Type | Data source | Parameters |
|---|---|---|---|
| ○ | S3 | s3://project1bysamir/csvfiles/ | Recrawl all |

▶ **Custom classifiers - *optional***
A classifier checks whether a given file is in a format the crawler can handle. If it is, the classifier creates a schema in the form of a StructType object that matches that data format.

[Cancel] [Previous] [Next]

---

✓ IAM Role "AWSGlueServiceRole-crawl" successfully updated
Successfully updated IAM Role "AWSGlueServiceRole-crawl". This role trusts AWS Glue and has permissions to access your AWS Glue Crawler targets.

# Configure security settings

## IAM role Info

Existing IAM role

AWSGlueServiceRole-crawl ▼ | C | [View ↗]

[Create new IAM role] [Update chosen IAM role]

Only IAM roles created by the AWS Glue console and have the prefix "AWSGlueServiceRole-" can be updated.

## Lake Formation configuration - *optional*
Allow the crawler to use Lake Formation credentials for crawling the data source. Learn more. ↗

☐ Use Lake Formation credentials for crawling S3 data source
Checking this box will allow the crawler to use Lake Formation credentials for crawling the data source. If the data source is registered in another account, you must provide the registered account ID. Otherwise, the crawler will crawl only those data sources associated to the account. Only applicable to S3, Glue Catalog, Iceberg, and Hudi data sources.

▶ **Security configuration - *optional***
Enable at-rest encryption with a security configuration.

[Cancel] [Previous] [Next]

---

## Databases (1)

Last updated (UTC)
March 17, 2024 at 22:54:22 | C | [Edit] [Delete] [Add database]

A database is a set of associated table definitions, organized into a logical group.

🔍 Filter databases

< 1 > ⚙

| | Name ▲ | Description ▽ | Location URI ▽ | Created on (UTC) ▽ |
|---|---|---|---|---|
| ☐ | crawl-output | - | - | March 17, 2024 at 22:54:19 |

# Set output and scheduling

## Output configuration  Info

Target database

crawl-output ▼    ⟳

[ Clear selection ]    [ Add database ⧉ ]

Table name prefix - *optional*

csvoutput-

Maximum table threshold - *optional*

This field sets the maximum number of tables the crawler is allowed to generate. In the event that this number is surpassed, the crawl will fail with an error. If not set, the crawler will automatically generate the number of tables depending on the data schema.

Type a number greater than 0

▶ Advanced options

## Crawler schedule

You can define a time-based schedule for your crawlers and jobs in AWS Glue. The definition of these schedules uses the Unix-like cron ⧉ syntax. Learn more ⧉.

Frequency

On demand ▼

[ Cancel ]    [ Previous ]    [ Next ]

---

**Step 1**
Set crawler properties

**Step 2**
Choose data sources and classifiers

**Step 3**
Configure security settings

**Step 4**
Set output and scheduling

**Step 5**
Review and create

# Review and create

## Step 1: Set crawler properties                    [ Edit ]

### Set crawler properties

| Name | Description | Tags |
|---|---|---|
| csv-crawler | crawl csv files in project1bysamir S3 bucket | - |

## Step 2: Choose data sources and classifiers       [ Edit ]

### Data sources (1)  Info

The list of data sources to be scanned by the crawler.

| Type | Data source | Parameters |
|---|---|---|
| S3 | s3://project1bysamir/csvfiles/data_utf... | Recrawl all |

## Step 3: Configure security settings                [ Edit ]

### Configure security settings

| IAM role | Security configuration | Lake Formation configuration |
|---|---|---|
| AWSGlueServiceRole-crawl | - | - |

---

## Step 3: Configure security settings                [ Edit ]

### Configure security settings

| IAM role | Security configuration | Lake Formation configuration |
|---|---|---|
| AWSGlueServiceRole-crawl | - | - |

## Step 4: Set output and scheduling                  [ Edit ]

### Set output and scheduling

| Database | Table prefix - *optional* | Maximum table threshold - *optional* | Schedule |
|---|---|---|---|
| crawl-output | csvoutput- | - | On demand |

[ Cancel ]    [ Previous ]    [ Create crawler ]

# csv-crawler

Run crawler   Edit   Delete

## Crawler properties

| | | | |
|---|---|---|---|
| **Name** | **IAM role** | **Database** | **State** |
| csv-crawler | AWSGlueServiceRole-crawl ↗ | crawl-output | READY |
| **Description** | **Security configuration** | **Lake Formation configuration** | **Table prefix** |
| crawl csv files in project1bysamir S3 bucket | - | - | csvoutput- |
| **Maximum table threshold** | | | |
| - | | | |

▶ Advanced settings

**Crawler runs** | Schedule | Data sources | Classifiers | Tags

### Crawler runs (1)
The list of crawler runs for this crawler.

Stop run   View CloudWatch logs ↗   View run details

🔍 Filter data          📅 Filter by a date and time range          ‹ 1 ›  ⚙

| | Start time (UTC) ▲ | End time (UTC) ▽ | Current/last duration ▽ | Status ▽ | DPU hours ▽ | Table changes ▽ |
|---|---|---|---|---|---|---|
| ○ | March 17, 2024 at 22:56:37 | March 17, 2024 at 22:57:35 | 57 s | ⊘ Completed | - | 1 table change, 0 partition changes |

---

AWS Glue > Tables > csvfiles

# csvfiles

Version 0 (Current version) ▼   Actions ▼

**Table overview** | Data quality New

**Table details** | Advanced properties

| | | | |
|---|---|---|---|
| **Name** | **Description** | **Database** | **Classification** |
| csvfiles | - | crawl-output | CSV |
| **Location** | **Connection** | **Deprecated** | **Last updated** |
| s3://project1bysamir/csvfiles/ | - | - | March 18, 2024 at 00:26:49 |

| | | |
|---|---|---|
| **Input format** | **Output format** | **Serde serialization lib** |
| org.apache.hadoop.mapred.TextInputFormat | org.apache.hadoop.hive.ql.io.HiveIgnoreKeyText OutputFormat | org.apache.hadoop.hive.serde2.lazy.LazySimpleS erDe |

**Schema** | Partitions | Indexes | Column statistics  - new

### Schema (8)
View and manage the table schema.

Edit schema as JSON   Edit schema

🔍 Filter schemas          ‹ 1 ›  ⚙

| # | Column name ▽ | Data type ▽ | Partition key ▽ | Comment ▽ |
|---|---|---|---|---|
| 1 | invoiceno | string | - | - |
| 2 | stockcode | string | - | - |
| 3 | description | string | - | - |
| 4 | quantity | bigint | - | - |
| 5 | invoicedate | string | - | - |
| 6 | unitprice | double | - | - |
| 7 | customerid | bigint | - | - |
| 8 | country | string | - | - |

## 3. Transform Data with AWS Glue Job:

Objective: Convert CSV data into Parquet format for efficient querying and analysis.

Action Items:

Create an AWS Glue job specifying the source data (the table generated by the crawler), the target format (Parquet), and the target location (a new or existing S3 bucket).

Choose or create an appropriate IAM role with the necessary permissions for the Glue job.

Write or generate the transformation script. AWS Glue can auto-generate a script for converting formats with minimal adjustments.

Run the Glue job to transform the CSV data into Parquet format.

# csv-to-parquet  ✎

⚠ Job has not been saved   Actions ▼   Save   Run

| Visual ① | Script | Job details | Runs | Data quality - *updated* | Schedules | Version Control |

+

Data source – Data Catalog
AWS Glue Data Catalog ✓

**Data source properties - Data Catalog**  ⛶

Name
AWS Glue Data Catalog

Database
Choose a database.
crawl-output ▼   ↻

▶ Use runtime parameters

Table
csvoutput-data_utf8_csv ▼   ↻

▶ Use runtime parameters

---

**Data preview**  |  Output schema  ▢ ▣

**Data preview** (200) Info  READY ⓘ   ↻   End session   Previewing 8 of 8 fields

🔍 Filter sample dataset  ⚙

| invoiceno ▽ | stockcode ▽ | description ▽ | quantity ▽ | invoicedate |
|---|---|---|---|---|
| | | ur | | |
| 536366 | 22632 | HAND WARMER RED POLKA DOT | 6 | 12/1/2010 8:2 |
| 536367 | 84879 | ASSORTED COLOUR BIRD ORNAMENT | 32 | 12/1/2010 8:3 |
| 536367 | 22745 | POPPY'S PLAYHOUSE BEDROOM | 6 | 12/1/2010 8:3 |
| | | POPPY'S PLAYHOUSE KIT | | |

---

# csv-to-parquet  ✎

⚠ Job has not been saved   Actions ▼   Save   Run

| **Visual** | Script | Job details | Runs | Data quality - *updated* | Schedules | Version Control |

+

Change Schema ✓

Data target - S3 bucket
Amazon S3 ✓

**Data target properties - S3**  ⛶

ApplyMapping - Transform

Format
Parquet ▼

ⓘ After you save your job, it will use Glue Studio's optimized Parquet writer.  ✕

Compression Type
Snappy ▼

S3 Target Location
Choose an S3 location in the format s3://bucket/prefix/object/ with a trailing slash (/).

🔍 s3://project1bysamir/parquetfolder/ ✕   View ⧉

Browse S3

Data Catalog update options   Info
Choose how you want to update the Data Catalog table's schema and partitions. These options will only apply if the Data Catalog table is an S3

---

**Data preview**  |  Output schema  ▢ ▣

ⓘ **Target node not supported**
You have selected a data target node which is not supported for data preview. Please select another type of node instead.

---

# csv-to-parquet

Last modified on 3/18/2024, 4:37:23 AM   Actions ▼   Save   **Run**

| Visual | Script | Job details | **Runs** | Data quality - *updated* | Schedules | Version Control |

**Job runs** (1/1) Info   Last updated (UTC) March 17, 2024 at 23:11:49   ↻   View details   Stop job run   **Table View** | Card View

🔍 Filter job runs by property   ‹ 1 ›  ⚙

| | Run status ▽ | Retries ▽ | Start time (UTC) ▽ | End time (UTC) ▽ | Duration ▽ | Capacity (DPUs) ▽ | Worker type ▽ | Glue version ▽ |
|---|---|---|---|---|---|---|---|---|
| ⦿ | ⊗ Failed | 0 | 2024/03/17 23:08:56 | 2024/03/17 23:11:21 | 2 m 13 s | 10 DPUs | G.1X | 4.0 |

**Run details** | Input arguments (10) | Continuous logs | Run insights | Metrics | Spark UI  ▢ ▣

⊗ Error Category: PERMISSION_ERROR; An error occurred while calling o125.pyWriteDynamicFrame. Access Denied (Service: Amazon S3; Status Code: 403; Error Code: AccessDenied; Request ID: DNAWS430PM DCVH3J; S3 Extended Request ID: iCixruVr+tCe+vgBfhk6VwY7Eo/Y5YE71yXHnMp7pUGvtzr63O0KuQhPEZUZ7yKrh3Sfx2bgFXM=; Proxy: null)

The Glue Job fails as Access was denied for the folder "Parquet" in the S3 bucket created in previous step. Add permissions, and the job is completed successfully.



✓ Policy AWSGlueServiceRole-crawl-EZCRC-s3Policy updated.                                                                              ×

Modify permissions in AWSGlueServiceRole-crawl-EZCRC-s3Policy Info
Add permissions by selecting services, actions, resources, and conditions. Build permission statements using the JSON editor.

Policy editor                                                    Visual    JSON    Actions ▾    ▣

```
1 ▾ {
2      "Version": "2012-10-17",
3      "Statement": [
4 ▾       {
5           "Effect": "Allow",
6 ▾         "Action": [
7               "s3:GetObject",
8               "s3:PutObject"
9           ],
10 ▾        "Resource": [
11              "arn:aws:s3:::project1bysamir/csvfiles/data_utf8.csv*",
12              "arn:aws:s3:::project1bysamir/parquetfolder/*"
13          ]
14        }
15     ]
16 }
```

Edit statement                          Remove

Add actions
Choose a service
🔍 Filter services

Included
S3

Available
AMP
API Gateway
API Gateway V2
ASC

csv-to-parquet                                          Last modified on 3/18/2024, 4:37:23 AM    Actions ▾    Save    Run

Visual    Script    Job details    Runs    Data quality - updated    Schedules    Version Control

Job runs (1/2) Info                                          Last updated (UTC)   ⟳   View details   Stop job run        Table View   Card View
                                                            March 17, 2024 at 23:16:26
🔍 Filter job runs by property                                                                                          < 1 >  ⚙

| | Run status | ▽ | Retries | ▽ | Start time (UTC) | ▽ | End time (UTC) | ▽ | Duration | ▽ | Capacity (DPUs) | ▽ | Worker type | ▽ | Glue version | ▽ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| ◉ | ✓ Succeeded | | 0 | | 2024/03/17 23:14:23 | | 2024/03/17 23:16:16 | | 1 m 38 s | | 10 DPUs | | G.1X | | 4.0 | |
| ○ | ✗ Failed | | 0 | | 2024/03/17 23:08:56 | | 2024/03/17 23:11:21 | | 2 m 13 s | | 10 DPUs | | G.1X | | 4.0 | |

aws    ▦ Services    🔍 Search                    [Alt+S]         ▣  △  ⑦  ⊙   Global ▾   Corestack_Role/swapna.samir.shukla.1988_gmail @ 0421-6739-7413 ▾

Amazon S3                ×      Amazon S3 > Buckets > project1bysamir > parquetfolder/

Buckets                         parquetfolder/                                                            📋 Copy S3 URI
Access Grants
Access Points                   Objects    Properties
Object Lambda Access Points
Multi-Region Access Points
Batch Operations                Objects (1) Info
IAM Access Analyzer for S3      ⟳   📋 Copy S3 URI   📋 Copy URL   ⬇ Download   Open ↗   Delete   Actions ▾   Create folder   📤 Upload
                                Objects are the fundamental entities stored in Amazon S3. You can use Amazon S3 inventory ↗ to get a list of all objects in your bucket. For others to access your objects, you'll need to explicitly
Block Public Access settings for   grant them permissions. Learn more ↗
this account
                                🔍 Find objects by prefix                                                           < 1 >  ⚙
▾ Storage Lens
Dashboards                      | □ | Name | ▲ | Type | ▽ | Last modified | ▽ | Size | ▽ | Storage class | ▽ |
Storage Lens groups             |---|---|---|---|---|---|---|---|---|---|---|
AWS Organizations settings      | □ | 📄 run-1710717332530-part-block-0-r-00000-snappy.parquet | | parquet | | March 18, 2024, 04:46:06 (UTC+05:30) | | 3.4 MB | | Standard | |

**4. Crawl Transformed Data:**

Objective: Catalog the metadata of the Parquet data files for querying.

Action Items:

Repeat the crawling process for the S3 location where the Parquet files are stored. This will allow AWS Athena to query the data efficiently.

Ensure the new crawler specifies the Parquet data's location, and after running, check that the Data Catalog contains the metadata for your Parquet files.

# Configure security settings

## IAM role  Info

Existing IAM role

[ AWSGlueServiceRole-crawl ▼ ]  [ C ]  [ View ⧉ ]

[ Create new IAM role ]  [ Update chosen IAM role ]

Only IAM roles created by the AWS Glue console and have the prefix "AWSGlueServiceRole-" can be updated.

## Lake Formation configuration - *optional*

Allow the crawler to use Lake Formation credentials for crawling the data source. Learn more. ⧉

☐ Use Lake Formation credentials for crawling S3 data source
Checking this box will allow the crawler to use Lake Formation credentials for crawling the data source. If the data source is registered in another account, you must provide the registered account ID. Otherwise, the crawler will crawl only those data sources associated to the account. Only applicable to S3, Glue Catalog, Iceberg, and Hudi data sources.

▶ Security configuration - *optional*
Enable at-rest encryption with a security configuration.

[ Cancel ]  [ Previous ]  [ Next ]

---

# Set output and scheduling

## Output configuration  Info

Target database

[ crawl-output-parquet ▼ ]  [ C ]

[ Clear selection ]  [ Add database ⧉ ]

Table name prefix - *optional*

[ parqout- ]

Maximum table threshold - *optional*
This field sets the maximum number of tables the crawler is allowed to generate. In the event that this number is surpassed, the crawl will fail with an error. If not set, the crawler will automatically generate the number of tables depending on the data schema.

[ Type a number greater than 0 ]

▶ Advanced options

## Crawler schedule

You can define a time-based schedule for your crawlers and jobs in AWS Glue. The definition of these schedules uses the Unix-like cron ⧉ syntax. Learn more ⧉.

Frequency

[ On demand ▼ ]

[ Cancel ]  [ Previous ]  [ Next ]

---

# parquet-crawler

Last updated (UTC)
March 17, 2024 at 23:22:33   [ C ]  [ Run crawler ]  [ Edit ]  [ Delete ]

## Crawler properties

| | | | |
|---|---|---|---|
| **Name**<br>parquet-crawler | **IAM role**<br>AWSGlueServiceRole-crawl ⧉ | **Database**<br>crawl-output-parquet | **State**<br>READY |
| **Description**<br>crawl the converted parquet file | **Security configuration**<br>- | **Lake Formation configuration**<br>- | **Table prefix**<br>parqout- |
| **Maximum table threshold**<br>- | | | |

▶ Advanced settings

| Crawler runs | Schedule | Data sources | Classifiers | Tags |

### Crawler runs (1)
The list of crawler runs for this crawler.

[ C ]  [ Stop run ]  [ View CloudWatch logs ⧉ ]  [ View run details ]

[ 🔍 Filter data ]   [ 📅 Filter by a date and time range ]                    < 1 > ⚙

| | Start time (UTC) ▲ | End time (UTC) ▽ | Current/last duration ▽ | Status ▽ | DPU hours ▽ | Table changes ▽ |
|---|---|---|---|---|---|---|
| ○ | March 17, 2024 at 23:22:37 | March 17, 2024 at 23:23:54 | 01 min 16 s | ⊘ Completed | - | 1 table change, 0 partition changes |

## 5. Analyze Data with AWS Athena:

Objective: Query the data to uncover insights such as the best-selling item and the geographical distribution of sales.

Action Items:

Open AWS Athena. Ensure it's configured to use the database generated by the Glue crawlers.

## Results (14)

| # | description | country |
|---|---|---|
| 1 | WORLD WAR 2 GLIDERS ASSTD DESIGNS | Spain |
| 2 | WORLD WAR 2 GLIDERS ASSTD DESIGNS | Germany |
| 3 | WORLD WAR 2 GLIDERS ASSTD DESIGNS | Switzerland |
| 4 | WORLD WAR 2 GLIDERS ASSTD DESIGNS | Sweden |
| 5 | WORLD WAR 2 GLIDERS ASSTD DESIGNS | France |
| 6 | WORLD WAR 2 GLIDERS ASSTD DESIGNS | Japan |
| 7 | WORLD WAR 2 GLIDERS ASSTD DESIGNS | Canada |
| 8 | WORLD WAR 2 GLIDERS ASSTD DESIGNS | Norway |
| 9 | WORLD WAR 2 GLIDERS ASSTD DESIGNS | United Kingdom |
| 10 | WORLD WAR 2 GLIDERS ASSTD DESIGNS | EIRE |
| 11 | WORLD WAR 2 GLIDERS ASSTD DESIGNS | Hong Kong |
| 12 | WORLD WAR 2 GLIDERS ASSTD DESIGNS | Portugal |
| 13 | WORLD WAR 2 GLIDERS ASSTD DESIGNS | Denmark |

| | A | B | C | D | E | F | G |
|---|---|---|---|---|---|---|---|
| 1 | description | country | | | | | |
| 2 | WORLD WAR 2 GLIDERS ASSTD DESIGNS | Spain | | | | | |
| 3 | WORLD WAR 2 GLIDERS ASSTD DESIGNS | Germany | | | | | |
| 4 | WORLD WAR 2 GLIDERS ASSTD DESIGNS | Switzerland | | | | | |
| 5 | WORLD WAR 2 GLIDERS ASSTD DESIGNS | Sweden | | | | | |
| 6 | WORLD WAR 2 GLIDERS ASSTD DESIGNS | France | | | | | |
| 7 | WORLD WAR 2 GLIDERS ASSTD DESIGNS | Japan | | | | | |
| 8 | WORLD WAR 2 GLIDERS ASSTD DESIGNS | Canada | | | | | |
| 9 | WORLD WAR 2 GLIDERS ASSTD DESIGNS | Norway | | | | | |
| 10 | WORLD WAR 2 GLIDERS ASSTD DESIGNS | United Kingdom | | | | | |
| 11 | WORLD WAR 2 GLIDERS ASSTD DESIGNS | EIRE | | | | | |
| 12 | WORLD WAR 2 GLIDERS ASSTD DESIGNS | Hong Kong | | | | | |
| 13 | WORLD WAR 2 GLIDERS ASSTD DESIGNS | Portugal | | | | | |
| 14 | WORLD WAR 2 GLIDERS ASSTD DESIGNS | Denmark | | | | | |
| 15 | WORLD WAR 2 GLIDERS ASSTD DESIGNS | Unspecified | | | | | |
| 16 | | | | | | | |
| 17 | | | | | | | |
| 18 | | | | | | | |
| 19 | | | | | | | |
| 20 | | | | | | | |
| 21 | | | | | | | |