```python
import pandas as pd
import numpy as np
import os

from cod_prep.claude.configurator import Configurator
from cod_prep.downloaders import (get_ages, get_cod_ages,
    get_current_location_hierarchy)
from cod_prep.utils.formatting import ages
from cod_prep.utils import report_if_merge_fail, print_log_message

CONF = Configurator()
SOURCE = "SOURCE_NAME"

L_DIR = "FILEPATH"

def read_in_data():

    df = pd.read_csv(L_DIR + 'FILENAME')

    df['cases'] = 1

    df['raw_specimen'] = df['source']
    df['raw_pathogen'] = df['species']

    df['sample_id'] = 'SOURCE-' + df['isolateid'].astype(str) + '-' + df['raw_specimen']

    df['year_id'] = df['year']

    return df


def format_age_group_id(df):

    age_dict = {
    '0 to 2 Years': 244,
    '3 to 12 Years': 291,
    '13 to 18 Years': 15,
    '19 to 64 Years': 163,
    '65 to 84 Years': 281,
    '85 and Over': 160,
    'Unknown': 283}

    df['age_group_id'] = df['agegroup'].map(age_dict)

    assert df.age_group_id.notnull().values.all()

    return df


def format_sex_id(df):
```

```python
52        df["sex_id"] = df["gender"].map({"Male": 1, "Female": 2, np.nan: 9})
53        report_if_merge_fail(df, "sex_id", "gender")
54
55        return df
56
57
58    def format_location_id(df):
59
60        location = get_current_location_hierarchy()
61
62        df['location_name'] = df['country']
63        df.loc[(df['location_name'] == 'United States') & (df['state'].notnull()),
      'location_name'] = df['state']
64
65        df = df.merge(location[['location_name','location_id']], how = 'left', on =
      'location_name')
66
67        df.loc[df['location_id'].isna(), 'location_name'].unique()
68
69        df.loc[df['location_name'] == 'United States', 'location_id'] = 102
70        df.loc[df['location_name'] == 'Hong Kong', 'location_id'] = 354
71        df.loc[df['location_name'] == 'Czech Republic', 'location_id'] = 47
72        df.loc[df['location_name'] == 'Russia', 'location_id'] = 62
73        df.loc[df['location_name'] == 'Venezuela', 'location_id'] = 133
74        df.loc[df['location_name'] == 'Korea, South', 'location_id'] = 68
75        df.loc[df['location_name'] == 'Taiwan', 'location_id'] = 8
76        df.loc[df['location_name'] == 'Vietnam', 'location_id'] = 20
77        df.loc[df['location_name'] == 'Slovak Republic', 'location_id'] = 54
78
79        assert df.location_id.notnull().values.all()
80
81        return df
82
83
84    def clean_drug(df):
85
86        antibio = df.iloc[:, 13:101]
87        antibio_names = [x for x in list(antibio) if "_i" in x]
88        antibio_numbers = [x for x in list(antibio) if (x not in list(antibio_names))]
89
90        df = df.drop(antibio_names, axis=1)
91
92        df.loc[df[antibio_numbers].isnull().apply(lambda x: all(x), axis=1), antibio_numbers] =
      'unknown'
93
94        keep = [x for x in list(df) if (x not in list(antibio_numbers))]
95
96        df = pd.melt(df, id_vars= keep, value_vars=antibio_numbers)
97
98        df = df[df['value'].notnull()]
99
100       df.rename(columns={'variable': 'raw_antibiotic', 'value': 'resistance'}, inplace=True)
101
102       assert df.raw_antibiotic.notnull().values.all()
```

```python
103        assert df.resistance.notnull().values.all()
104
105        return df
106
107
108  def clean_SOURCE():
109        df = read_in_data()
110        df = format_age_group_id(df)
111        df = format_sex_id(df)
112        df = format_location_id(df)
113        df = clean_drug(df)
114
115        df["nid"] = 410524
116
117        return df
118
119
120  if __name__ == '__main__':
121        df = clean_SOURCE()
122        demo_cols = ['nid', 'sample_id', 'location_id', 'year_id', 'age_group_id', 'sex_id']
123        biology = ['raw_pathogen', 'raw_antibiotic', 'resistance', 'raw_specimen']
124        other_values = ['cases']
125        lowercase_cols = ['raw_pathogen', 'raw_antibiotic', 'raw_specimen']
126        df = df[demo_cols + biology + other_values]
127
128        # Lower case a few mapped columns
129        for col in lowercase_cols:
130            df[col] = df[col].str.lower().str.strip()
131            assert df[col].notnull().values.all()
132
133        df.to_csv("FILEPATH", index=False)
```