# COMP101 — Week 10

# File Handling in Python

Reading, writing, and cleaning external data

UM6P — SASE

December 19, 2025

- Open, read, and write text files safely in Python.

- Open, read, and write text files safely in Python.
- Use context managers (with open(...) as f) to avoid resource leaks.

- Open, read, and write text files safely in Python.
- Use context managers (with open(...) as f) to avoid resource leaks.
- Process line-based data: strip, split, convert types, validate.

# Today's Goals

- Open, read, and write text files safely in Python.
- Use context managers (with open(...) as f) to avoid resource leaks.
- Process line-based data: strip, split, convert types, validate.
- Handle common file errors (FileNotFoundError, malformed lines).

- Open, read, and write text files safely in Python.
- Use context managers (with open(...) as f) to avoid resource leaks.
- Process line-based data: strip, split, convert types, validate.
- Handle common file errors (FileNotFoundError, malformed lines).
- Think about performance: reading whole files vs streaming line by line.

- Programs rarely work only with data you type at the keyboard.

- Programs rarely work only with data you type at the keyboard.
- Real tasks: logs, configuration files, CSV exports, reports.

- Programs rarely work only with data you type at the keyboard.
- Real tasks: logs, configuration files, CSV exports, reports.
- Files let us:

- Programs rarely work only with data you type at the keyboard.
- Real tasks: logs, configuration files, CSV exports, reports.
- Files let us:
  - Persist results between program runs.

- Programs rarely work only with data you type at the keyboard.
- Real tasks: logs, configuration files, CSV exports, reports.
- Files let us:
  - Persist results between program runs.
  - Reuse existing datasets.

- Programs rarely work only with data you type at the keyboard.
- Real tasks: logs, configuration files, CSV exports, reports.
- Files let us:
  - Persist results between program runs.
  - Reuse existing datasets.
  - Exchange data with other tools (Excel, databases, etc.).

- Programs rarely work only with data you type at the keyboard.
- Real tasks: logs, configuration files, CSV exports, reports.
- Files let us:
  - Persist results between program runs.
  - Reuse existing datasets.
  - Exchange data with other tools (Excel, databases, etc.).
- Input/output is the first point where your code meets the messy outside world.

### Definition

A file is a named sequence of bytes stored by the operating system on a device (disk, SSD, USB, . . . ). Python gives you a higher-level interface to read and write these bytes.

- The OS controls access to files (permissions, paths, modes).

### Definition

A file is a named sequence of bytes stored by the operating system on a device (disk, SSD, USB, ...). Python gives you a higher-level interface to read and write these bytes.

- The OS controls access to files (permissions, paths, modes).
- Python wraps this with a file object.

#### Definition

A file is a named sequence of bytes stored by the operating system on a device (disk, SSD, USB, . . . ). Python gives you a higher-level interface to read and write these bytes.

- The OS controls access to files (permissions, paths, modes).
- Python wraps this with a file object.
- For text files, Python converts between bytes and strings using an encoding.

- Basic usage:

```
f = open("data.txt", "r")   # "r" = read text
text = f.read()
f.close()
```

- Basic usage:

```
f = open("data.txt", "r")    # "r" = read text
text = f.read()
f.close()
```

- Filename: string, relative or absolute path.

- Basic usage:

```python
f = open("data.txt", "r")    # "r" = read text
text = f.read()
f.close()
```

- Filename: string, relative or absolute path.
- Mode (common):

- Basic usage:

```
f = open("data.txt", "r")    # "r" = read text
text = f.read()
f.close()
```

- Filename: string, relative or absolute path.
- Mode (common):
    - "r": read (file must exist)

- Basic usage:

```
f = open("data.txt", "r")    # "r" = read text
text = f.read()
f.close()
```

- Filename: string, relative or absolute path.
- Mode (common):
    - "r": read (file must exist)
    - "w": write (overwrite if file exists)

- Basic usage:

```
f = open("data.txt", "r")    # "r" = read text
text = f.read()
f.close()
```

- Filename: string, relative or absolute path.
- Mode (common):
  - "r": read (file must exist)
  - "w": write (overwrite if file exists)
  - "a": append (add to end)

- Basic usage:

```
f = open("data.txt", "r")    # "r" = read text
text = f.read()
f.close()
```

- Filename: string, relative or absolute path.
- Mode (common):
    - "r": read (file must exist)
    - "w": write (overwrite if file exists)
    - "a": append (add to end)
    - "rb"/"wb": binary modes

- Basic usage:

```
f = open("data.txt", "r")    # "r" = read text
text = f.read()
f.close()
```

- Filename: string, relative or absolute path.
- Mode (common):
    - "r": read (file must exist)
    - "w": write (overwrite if file exists)
    - "a": append (add to end)
    - "rb"/"wb": binary modes
- Calling close() releases the resource; forgetting it is a bug.

#### Definition

A context manager automatically acquires and releases a resource, even if an error occurs inside the block.

```
# Recommended pattern
with open("data.txt", "r", encoding="utf-8") as f:
    text = f.read()
    # work with text here

# f is automatically closed here
```

- Always prefer with for files.

### Definition

A context manager automatically acquires and releases a resource, even if an error occurs inside the block.

```python
# Recommended pattern
with open("data.txt", "r", encoding="utf-8") as f:
    text = f.read()
    # work with text here

# f is automatically closed here
```

- Always prefer with for files.
- Avoids leaking open file handles.

### Definition

A context manager automatically acquires and releases a resource, even if an error occurs inside the block.

```python
# Recommended pattern
with open("data.txt", "r", encoding="utf-8") as f:
    text = f.read()
    # work with text here

# f is automatically closed here
```

- **Always prefer** with for files.
- Avoids leaking open file handles.
- Makes the lifetime of the file object explicit.

- Whole file at once (read()):

```python
with open("data.txt", "r", encoding="utf-8") as f:
    data = f.read()    # one big string
```

```python
with open("data.txt", "r", encoding="utf-8") as f:
    for line in f:           # streaming iteration
        line = line.strip()
        # process line
```

```python
with open("data.txt", "r", encoding="utf-8") as f:
    lines = f.readlines()   # list of strings
```

- Whole file at once (read()):

```python
with open("data.txt", "r", encoding="utf-8") as f:
    data = f.read()   # one big string
```

- Line by line (readline() or iteration):

```python
with open("data.txt", "r", encoding="utf-8") as f:
    for line in f:          # streaming iteration
        line = line.strip()
        # process line
```

```python
with open("data.txt", "r", encoding="utf-8") as f:
    lines = f.readlines()  # list of strings
```

- Whole file at once (read()):

```python
with open("data.txt", "r", encoding="utf-8") as f:
    data = f.read()   # one big string
```

- Line by line (readline() or iteration):

```python
with open("data.txt", "r", encoding="utf-8") as f:
    for line in f:          # streaming iteration
        line = line.strip()
        # process line
```

- Read all lines into a list:

```python
with open("data.txt", "r", encoding="utf-8") as f:
    lines = f.readlines()  # list of strings
```

# Writing To Files

```python
# Overwrite or create file
with open("results.txt", "w", encoding="utf-8") as f:
    f.write("Average score: 15.7\n")
    f.write("Median score: 14.0\n")

# Append to existing file
with open("results.txt", "a", encoding="utf-8") as f:
    f.write("New run finished.\n")
```

- "w" truncates the file if it exists.

# Writing To Files

```python
# Overwrite or create file
with open("results.txt", "w", encoding="utf-8") as f:
    f.write("Average score: 15.7\n")
    f.write("Median score: 14.0\n")

# Append to existing file
with open("results.txt", "a", encoding="utf-8") as f:
    f.write("New run finished.\n")
```

- "w" truncates the file if it exists.
- "a" keeps existing content and adds at the end.

```python
# Overwrite or create file
with open("results.txt", "w", encoding="utf-8") as f:
    f.write("Average score: 15.7\n")
    f.write("Median score: 14.0\n")

# Append to existing file
with open("results.txt", "a", encoding="utf-8") as f:
    f.write("New run finished.\n")
```

- "w" truncates the file if it exists.
- "a" keeps existing content and adds at the end.
- Always control your newlines "\n" explicitly.

- Text in Python is Unicode; files are bytes.

```python
with open("data.txt", "r", encoding="utf-8") as f:
    text = f.read()
```

- Text in Python is Unicode; files are bytes.
- Encoding tells Python how to translate between them:

```python
with open("data.txt", "r", encoding="utf-8") as f:
    text = f.read()
```

- Text in Python is Unicode; files are bytes.

- Encoding tells Python how to translate between them:

```python
with open("data.txt", "r", encoding="utf-8") as f:
    text = f.read()
```

- Use encoding="utf-8" unless you have a strong reason not to.

- Text in Python is Unicode; files are bytes.

- Encoding tells Python how to translate between them:

```python
with open("data.txt", "r", encoding="utf-8") as f:
    text = f.read()
```

- Use encoding="utf-8" unless you have a strong reason not to.

- Different systems use different newline conventions; Python normalizes them to "\n" when reading in text mode.

- Text in Python is Unicode; files are bytes.

- Encoding tells Python how to translate between them:

```python
with open("data.txt", "r", encoding="utf-8") as f:
    text = f.read()
```

- Use encoding="utf-8" unless you have a strong reason not to.

- Different systems use different newline conventions; Python normalizes them to "\n" when reading in text mode.

- Encoding mistakes show up as strange symbols or UnicodeDecodeError.

```python
with open("scores.txt", "r", encoding="utf-8") as f:
    scores = []
    for line in f:
        line = line.strip()
        if not line:           # skip empty lines
            continue
        parts = line.split(",")   # e.g. "Ali,17"
        name = parts[0]
        score = int(parts[1])
        scores.append((name, score))
```

- strip() to remove extra spaces and newlines.

```python
with open("scores.txt", "r", encoding="utf-8") as f:
    scores = []
    for line in f:
        line = line.strip()
        if not line:          # skip empty lines
            continue
        parts = line.split(",")    # e.g. "Ali,17"
        name = parts[0]
        score = int(parts[1])
        scores.append((name, score))
```

- strip() to remove extra spaces and newlines.
- split() to break structured text into fields.

```python
with open("scores.txt", "r", encoding="utf-8") as f:
    scores = []
    for line in f:
        line = line.strip()
        if not line:            # skip empty lines
            continue
        parts = line.split(",")     # e.g. "Ali,17"
        name = parts[0]
        score = int(parts[1])
        scores.append((name, score))
```

- strip() to remove extra spaces and newlines.
- split() to break structured text into fields.
- Type conversion (e.g. int(), float()) with basic validation.

```python
filename = "scores.txt"

try:
    with open(filename, "r", encoding="utf-8") as f:
        data = f.read()
except FileNotFoundError:
    print(f"File {filename} not found.")
except PermissionError:
    print(f"No permission to read {filename}.")
```

- Expect missing files and permission issues.

# Handling File Errors

```python
filename = "scores.txt"

try:
    with open(filename, "r", encoding="utf-8") as f:
        data = f.read()
except FileNotFoundError:
    print(f"File {filename} not found.")
except PermissionError:
    print(f"No permission to read {filename}.")
```

- Expect missing files and permission issues.
- Catch specific exceptions, not a generic except:.

```python
filename = "scores.txt"

try:
    with open(filename, "r", encoding="utf-8") as f:
        data = f.read()
except FileNotFoundError:
    print(f"File {filename} not found.")
except PermissionError:
    print(f"No permission to read {filename}.")
```

- Expect missing files and permission issues.
- Catch specific exceptions, not a generic except:.
- Decide: fail with a clear message, or create a default file.

### Algorithm

Task: Read a file scores.txt with lines "name,score", ignore invalid lines, and write a cleaned file plus summary.

```python
valid = []
with open("scores.txt", "r", encoding="utf-8") as f:
    for line in f:
        line = line.strip()
        if not line:
            continue
        parts = line.split(",")
        if len(parts) != 2:
            continue
        name, raw = parts
        try:
            score = float(raw)
        except ValueError:
            continue
        valid.append((name, score))
```

```python
with open("scores_clean.txt", "w", encoding="utf-8") as out:
    for name, score in valid:
        out.write(f"{name},{score}\n")
```

- Small files: read() or readlines() are fine.

```
with open("big_log.txt", "r", encoding="utf-8") as f:
    for line in f:
        process(line)
```

- Small files: read() or readlines() are fine.
- Large files: prefer iterating line by line:

```python
with open("big_log.txt", "r", encoding="utf-8") as f:
    for line in f:
        process(line)
```

- Small files: read() or readlines() are fine.

- Large files: prefer iterating line by line:

```
with open("big_log.txt", "r", encoding="utf-8") as f:
    for line in f:
        process(line)
```

- This keeps memory usage almost constant.

- Small files: read() or readlines() are fine.

- Large files: prefer iterating line by line:

```python
with open("big_log.txt", "r", encoding="utf-8") as f:
    for line in f:
        process(line)
```

- This keeps memory usage almost constant.

- Think about complexity: your algorithm may be fast, but IO can dominate.

- Forgetting to close files (fixed by using with).

- Forgetting to close files (fixed by using with).
- Assuming the file always exists and is well-formed.

- Forgetting to close files (fixed by using with).
- Assuming the file always exists and is well-formed.
- Ignoring leading/trailing spaces and blank lines.

- Forgetting to close files (fixed by using with).
- Assuming the file always exists and is well-formed.
- Ignoring leading/trailing spaces and blank lines.
- Confusing text and binary modes.

- Forgetting to close files (fixed by using with).
- Assuming the file always exists and is well-formed.
- Ignoring leading/trailing spaces and blank lines.
- Confusing text and binary modes.
- Accidentally overwriting files with "w" instead of "a".

- Forgetting to close files (fixed by using with).
- Assuming the file always exists and is well-formed.
- Ignoring leading/trailing spaces and blank lines.
- Confusing text and binary modes.
- Accidentally overwriting files with "w" instead of "a".
- Hard-coding absolute paths that only work on your machine.

# Summary

- File handling connects your code to persistent data.

# Summary

- File handling connects your code to persistent data.
- Use with open(..., mode, encoding) as your default pattern.

- **File handling** connects your code to persistent data.
- Use with open(..., mode, encoding) as your default pattern.
- Read and write text carefully: strip, split, convert, validate.

- File handling connects your code to persistent data.
- Use with open(..., mode, encoding) as your default pattern.
- Read and write text carefully: strip, split, convert, validate.
- Handle file-related exceptions explicitly.

- **File handling** connects your code to persistent data.
- Use **with open(..., mode, encoding)** as your default pattern.
- Read and write text carefully: strip, split, convert, validate.
- Handle file-related exceptions explicitly.
- For larger files, **stream line by line** instead of loading everything.