

# Gradio for ML Models

TimeSeries Team

Ngày 22 tháng 9 năm 2025

## Mục lục

<b>1</b>	<b>Giới thiệu &amp; Mục tiêu</b>	<b>3</b>
<b>2</b>	<b>Vì sao dùng Gradio?</b>	<b>6</b>
<b>3</b>	<b>Khung tư duy: Input – Model – Output</b>	<b>9</b>
3.1	Đặt “hộp đồng dữ liệu” ngay từ đầu	9
3.2	Ánh xạ Input ↔ Component	9
3.3	Thiết kế hàm <code>model_fn</code> : chữ ký, kiểm tra, trả về	9
3.4	Tiền xử lý (preprocess) & Hậu xử lý (postprocess)	10
3.5	Thiết kế Output: dễ đọc, kiểm thử, và nhất quán	10
3.6	Ví dụ 1: NLP đơn giản (text + slider → text)	11
3.7	Ví dụ 2: Ảnh (image + threshold → image + text)	11
3.8	Nhiều Input/Output: quy ước <b>thứ tự</b> và <b>kiểu</b>	12
3.9	Xử lý lỗi và trạng thái biên (edge cases)	12
3.10	Hiệu năng và ổn định	12
3.11	Bảo mật dữ liệu (ở mức demo)	12
<b>4</b>	<b>Các component hay dùng</b>	<b>14</b>
4.1	Nguyên tắc chọn component	14
4.2	Nhóm Input cơ bản	14
4.3	Nhóm Output phổ biến	15
4.4	Mẫu ánh xạ nhanh Input → Output	15
4.5	Ví dụ 1: Textbox + Slider → Textbox + JSON	16
4.6	Ví dụ 2: Image + Slider → Image + Textbox	16
4.7	Xử lý file và media an toàn	17
4.8	Mẹo UI/UX nhỏ nhưng hiệu quả	17
<b>5</b>	<b>Hai cách dựng giao diện: Interface vs. Blocks</b>	<b>18</b>
5.1	Interface: đường tắt cho MVP/POC	18
5.2	Blocks: “lego” để mở rộng UI/UX và luồng	19
5.3	Tiêu chí chọn nhanh	20
5.4	Mẫu nâng cấp: từ Interface → Blocks	20
5.5	Các mẫu thiết kế (patterns) thường gặp	20
5.6	Bẫy thường gặp và cách tránh	21
5.7	Hiệu năng & khả năng mở rộng	21
<b>6</b>	<b>Bố cục UI: Row/Column, Tab, Accordion</b>	<b>22</b>
6.1	Nguyên tắc bố cục tổng quát	22
6.2	Row/Column: xương sống của layout	22
6.3	Tab: phân tách chức năng theo ngữ nghĩa	23
6.4	Accordion: giấu bớt phần ít dùng	23
6.5	Tổ chức “sidebar” tham số	24

6.6	Chia nhóm “kết quả chính” và “siêu dữ liệu”	24
6.7	Anti-patterns cần tránh	24
6.8	Mẹo vi mô để UI gọn gàng	25
<b>7</b>	<b>Event listeners: chọn đúng để app mượt</b>	<b>26</b>
7.1	Bốn sự kiện cốt lõi	26
7.2	Quy tắc chọn nhanh	26
7.3	Mẫu kết nối sự kiện chuẩn	26
7.4	Tránh lặp gọi với <b>Slider</b>	27
7.5	Ghép nhiều sự kiện cho một luồng	27
7.6	Tách <i>preprocess/infer/postprocess</i> bằng sự kiện	27
7.7	Xử lý lỗi & trạng thái biên qua sự kiện	28
7.8	Hàng đợi ( <b>queue=True</b> ) và phản hồi người dùng	28
7.9	Anti-patterns thường gặp	28
<b>8</b>	<b>Thực hành: Object Detection (YOLOv8n) với Gradio</b>	<b>29</b>
8.1	Chuẩn bị môi trường	29
8.2	Cấu trúc dự án gợi ý	29
8.3	Hằng số và tiện ích xử lý ảnh	29
8.4	Nạp mô hình một lần ở phạm vi toàn cục	30
8.5	Hàm suy luận: annotate ảnh + tóm tắt kết quả	30
8.6	Hàm tiện ích: Clear và tải ảnh mẫu	31
8.7	Dựng UI bằng <b>Blocks</b>	31
8.8	Kịch bản kiểm thử nhanh	32
8.9	Tối ưu hiệu năng	32
8.10	Xử lý lỗi và tình huống biên	32
8.11	An toàn dữ liệu (mức demo)	32
8.12	Checklist hoàn tất phần thực hành	33
<b>9</b>	<b>Triển khai lên Hugging Face Spaces</b>	<b>34</b>
9.1	Kiến trúc nền tảng nhanh	34
9.2	Chuẩn bị dự án	34
9.3	Tạo Space và đẩy mã	35
9.4	Cấu hình nâng cao	35
9.5	Khắc phục sự cố thường gặp	35
9.6	Checklist triển khai thành công	35
<b>10</b>	<b>Kết luận</b>	<b>36</b>

# 1 Giới thiệu & Mục tiêu

## Bức tranh tổng quan

Bạn có một mô hình hoặc đơn giản là một hàm suy luận (inference function) trong Python và muốn **cho người khác dùng thử ngay** mà không cần tự xây toàn bộ frontend/backend. Mục tiêu của bài viết là giúp bạn đưa mô hình đó thành **ứng dụng web demo** qua Gradio theo lộ trình rõ ràng, có thể tái sử dụng cho nhiều bài toán khác nhau (thị giác máy tính, NLP, tabular...).

## Tư duy cốt lõi: Input – Model – Output

Trong ngữ cảnh suy luận, ta gom bài toán về ba mảnh ghép:

**Input** Người dùng sẽ cung cấp dữ liệu gì? (văn bản/số/ảnh/tệp/lựa chọn...). Xác định đúng giúp ta **chọn component Gradio** phù hợp (vd. Textbox, Slider, Image).

**Model** Hàm Python nhận input và trả về kết quả. Hãy **định dạng vào/ra** ngay từ đầu (PIL Image, NumPy array, str, dict...) để tránh chuyển đổi phức tạp về sau.

**Output** Ta muốn hiển thị gì cho người dùng? (label/text, ảnh đã vẽ, bảng, JSON, tệp tải về...). Với kết quả chỉ để xem, nên đặt `interactive=false` để tránh chĩnh nhảm.

Cách tư duy này giúp bạn tách bạch phần UI (component) khỏi phần logic (hàm suy luận), từ đó **giảm thời gian ghép nối và dễ mở rộng**.

## Vì sao chọn Gradio cho demo nhanh?

- **Một ngôn ngữ (Python) cho cả UI lẫn model**: không cần HTML/CSS/JS hoặc framework web riêng.
- **Component dựng sẵn**: Textbox, Slider, Image, File, Dataframe, JSON...; có thể dùng cho cả input & output.
- **Sự kiện rõ ràng**: `.click()`, `.change()`, `.submit()`, `.release()` giúp điều khiển thời điểm gọi hàm chính xác.
- **Chia sẻ nhanh**: `launch(share=true)` tạo liên kết dùng tạm thời; có thể nâng cấp lên Hugging Face Spaces để có URL bền vững.

## Bạn sẽ xây được gì sau bài này?

1. Hiểu và áp dụng khung **Input – Model – Output** để thiết kế giao diện tối thiểu mà hiệu quả.
2. Dựng demo với **Interface** (một hàm duy nhất) và nâng cấp lên **Blocks** (bố cục linh hoạt, nhiều sự kiện).
3. Tổ chức layout bằng **Row/Column**, thêm các nút **Sample/Clear**, và gán sự kiện hợp lý (ưu tiên `.release` cho slider khi tác vụ nặng).
4. Xây một **demo Object Detection** gọn nhẹ (YOLOv8n) và hiển thị ảnh *đã annotate* kèm tóm tắt kết quả.
5. Đóng gói và **đưa demo “lên sóng”** bằng Hugging Face Spaces.

## Đối tượng & yêu cầu tiên quyết

- Đọc giả đã quen với Python cơ bản (hàm, `virtualenv`/`conda`), biết chạy notebook hoặc script.
- Có tối thiểu kiến thức nhập môn về mô hình ML đang sử dụng (cách nạp model, chạy suy luận).
- Môi trường đề xuất: Python 3.10+, `pip` cập nhật; GPU là *tùy chọn* (demo dùng bản YOLO nhỏ có thể chạy CPU).

## Phạm vi & không nằm trong phạm vi

**Phạm vi** Dựng **demo suy luận** chạy được, có giao diện; không nhằm tối ưu training hay MLOps phức tạp.

**Không phạm vi** Không đi sâu huấn luyện, không chuẩn hóa CI/CD hay autoscaling; không bàn chi tiết tối ưu kiến trúc mạng.

## Kiến trúc tinh gọn (tư duy hệ thống)

**Mục tiêu:** thay vì xây “FE ↔ API ↔ Model” đầy đủ, ta *gộp* vào một tiến trình Python:

1. Nạp trọng số/mô hình.
2. Định nghĩa hàm suy luận chuẩn hóa *đầu vào/đầu ra*.
3. Khai báo UI bằng Gradio; kết nối sự kiện tới hàm suy luận.
4. Khởi chạy server cục bộ và/hoặc chia sẻ liên kết.

Cách làm này giúp **rút ngắn thời gian** từ notebook tới ứng dụng có thể tương tác, đủ tốt cho việc truyền thông nội bộ, trình bày, thu thập phản hồi.

## Lộ trình nội dung của bài

1. Vì sao Gradio: lợi ích thực tế khi làm demo nhanh.
2. Khung **Input – Model – Output** và cách ánh xạ sang component.
3. **Interface** vs **Blocks**: khi nào dùng cái nào; ví dụ tối thiểu.
4. Bố cục với **Row/Column**, bổ sung **Tab/Accordion** khi cần.
5. Event listeners và các lựa chọn: `.click`, `.change`, `.submit`, `.release`.
6. Thực hành: demo Object Detection (YOLOv8n) + Gradio.
7. Đưa demo lên Hugging Face Spaces.

## Tiêu chí đánh giá một bản demo “đủ tốt”

- **Rõ ràng:** hiển thị Input & Output mạch lạc; có mô tả ngắn cách dùng.
- **Ổn định:** thao tác phổ biến (upload, kéo slider, bấm nút) không gây lỗi; có nút **Clear**.
- **Nhẹ:** chạy được trên CPU/GPU phổ thông; nếu ảnh lớn, có bước resize.
- **Chia sẻ được:** `share=true` hoặc deploy lên Spaces; ghi lại `requirements.txt`.

## “Hello, Gradio!” tối thiểu (tùy chọn chạy thử)

Đoạn mã cực ngắn dưới đây giúp bạn kiểm tra nhanh môi trường trước khi đi sâu:

```
1 import gradio as gr
2
3 def echo(x: str) -> str:
4     return f"Bạn vừa nhập: {x}"
5
6 demo = gr.Interface(
7     fn=echo,
8     inputs=gr.Textbox(label="Nhập gì đó"),
9     outputs=gr.Textbox(label="Phản hồi", interactive=False),
10    title="Hello Gradio",
11    description="Demo tối thiểu để kiểm tra môi trường."
12 )
13
14 if __name__ == "__main__":
15     demo.launch(share=False) # chuyển sang True nếu muốn có link tạm công khai
```

Listing 1: Hello Gradio: kiểm tra môi trường & cách launch tối thiểu

## Checklist thiết lập nhanh trước khi thực hành

- Python  $\geq 3.10$ , pip cập nhật: `python -V`, `pip -V`.
- Tạo môi trường ảo (khuyến nghị) và `pip install gradio`.
- Chạy “Hello Gradio” ở trên để xác nhận giao diện mở được.
- Với bài thực hành Object Detection: cài thêm `ultralytics`, `opencv-python`, `pillow`, `torch`, `torchvision`.

## Kết nối sang các phần tiếp theo

Sau phần mở đầu, ta sẽ lần lượt:

1. Chọn `Interface` hay `Blocks` cho bài toán cụ thể.
2. Bố cục UI mạch lạc với `Row/Column` và các tiện ích.
3. Gắn sự kiện hợp lý để tối ưu trải nghiệm và tài nguyên.
4. Áp dụng tất cả vào demo Object Detection và triển khai “lên sóng”.

## 2 Vì sao dùng Gradio?

### Bài toán thực tế khi đưa mô hình ra dùng thử

Khi có một mô hình hoặc hàm suy luận, bạn thường muốn:

1. **Cho người khác dùng ngay** (đồng đội, giảng viên, khách hàng) mà không cần cài môi trường.
2. **Thu thập phản hồi sớm**: dữ liệu đầu vào đa dạng, góc nhìn người dùng thật, ghi nhận vấn đề và cải tiến.
3. **Trình diễn nhanh** trong workshop, tiết giảng, hoặc buổi demo nội bộ.

Nếu đi theo hướng truyền thống (FE-BE-Model), thời gian dựng hạ tầng sẽ lấn át thời gian tinh chỉnh mô hình. Gradio giải quyết nút thắt này bằng cách **gom UI + server vào Python**.

### Lợi ích cốt lõi của Gradio

- **Một ngôn ngữ duy nhất (Python)**: bạn vừa gọi model, vừa dựng UI, vừa chạy server — tất cả trong một script.
- **Component UI dựng sẵn**: Textbox, Slider, Image, File, Dataframe, JSON, ... giúp mô tả *đúng* kiểu dữ liệu.
- **Event rõ ràng**: `.click()`, `.change()`, `.submit()`, `.release()` cho phép điều khiển thời điểm chạy suy luận chính xác, tránh gọi thừa.
- **Chia sẻ tức thì**: `launch(share=True)` tạo liên kết tạm công khai; dễ *mời người khác dùng thử* không cần cài đặt.
- **Triển khai tối giản**: có thể đẩy lên Hugging Face Spaces để có URL bền vững, đóng vai trò như “sandbox” để demo/POC.

### So sánh nhanh với các hướng tiếp cận khác

Tiêu chí	Gradio	FE-BE tự xây (VD: React + FastAPI)
Thời gian lên demo	<b>Rất nhanh</b> (vài chục dòng)	Chậm hơn (tách repo FE, BE, CORS, build ...)
Độ linh hoạt UI/UX	Vừa đủ cho demo, POC	<b>Rất cao</b> (thiết kế tùy ý, scale lớn)
Độ phức tạp vận hành	Thấp (một tiến trình)	Cao (deploy nhiều dịch vụ, CI/CD)
Kiểm soát bảo mật	Vừa phải (link tạm, Spaces)	<b>Toàn quyền</b> (auth, gateway, secrets)
Thích hợp cho	Học tập, demo nhanh, nội bộ	Sản phẩm production, SLA, traffic lớn

### Khi nào nên dùng Gradio?

- **Demo nhanh/MVP**: một hàm suy luận, một vài tham số, cần link để người khác thử.
- **Bài giảng/workshop**: tập trung vào mô hình hơn là hạ tầng, sinh viên chạy được ngay.
- **Nội bộ/POC**: xác nhận tính hữu dụng trước khi đầu tư xây hạ tầng bài bản.

## Khi nào chưa phù hợp?

- **Yêu cầu UI/UX phức tạp:** workflow nhiều bước, realtime đa thiết bị, brand guideline chặt.
- **Ràng buộc bảo mật/ngân hàng/y tế:** cần chuẩn hoá IAM, audit log, phân vùng mạng nghiêm ngặt.
- **SLA/traffic lớn:** cần auto-scaling, canary, quan sát hệ thống đầy đủ (metrics, tracing).

## Mô hình kiến trúc tối giản (tư duy hệ thống)

Gradio hướng tới **đường đi ngắn nhất** từ mã Python tới app có thể tương tác:

1. **Nạp mô hình/trọng số** khi khởi động tiến trình.
2. **Định nghĩa hàm suy luận** (chuẩn hoá vào/ra, xử lý lỗi, giới hạn kích thước dữ liệu).
3. **Khai báo UI:** chọn component khớp kiểu dữ liệu; tạo bố cục bằng Row/Column/Tab.
4. **Gắn sự kiện:** kết nối component → hàm suy luận bằng `.click()`, `.release()`, ...
5. **Launch:** chạy server cục bộ; tùy chọn `share=true` để chia sẻ nhanh.

## Ví dụ so sánh độ phức tạp (ý tưởng)

Không cần *routing*, *schema request/response*, *template HTML*. Thay vào đó:

```

1 import gradio as gr
2
3 def predict(text, temp: float):
4     # ... gọi model NLP/LLM/Classifier ...
5     return f"Output voi temp={temp:.2f}: {text[::-1]}"
6
7 with gr.Blocks(title="Demo nhanh") as demo:
8     with gr.Row():
9         with gr.Column(scale=1):
10             inp = gr.Textbox(label="Nhập văn bản")
11             temp = gr.Slider(0.0, 1.0, value=0.5, step=0.05, label="Temperature")
12             run = gr.Button("Chạy")
13         with gr.Column(scale=2):
14             out = gr.Textbox(label="Kết quả", interactive=False, lines=6)
15     run.click(predict, inputs=[inp, temp], outputs=out)
16
17 demo.launch(share=False)

```

Listing 2: Từ hàm Python đến app có giao diện trong vài dòng

## Hiệu năng & trải nghiệm người dùng (UX)

- **Giảm số lần suy luận thừa:** dùng `.release` với Slider thay vì `.change` cho tác vụ nặng.
- **Hạn chế kích thước dữ liệu:** resize ảnh lớn trước khi suy luận; từ chối tệp quá dung lượng.
- **Phản hồi rõ ràng:** dùng `interactive=false` cho output; hiển thị thông điệp “đang xử lý” khi tác vụ lâu.
- **Hàng đợi:** bật `queue=True` nếu dự kiến nhiều yêu cầu đồng thời.

## Bảo mật & chia sẻ

- **share=true**: liên kết tạm thời, tiện thử nghiệm; không nên dùng cho dữ liệu nhạy cảm.
- **Hugging Face Spaces**: URL công khai bền vững; có thể cấu hình *secrets* và quyền truy cập.
- **Dọn dữ liệu tạm**: đảm bảo không lưu trữ bản sao dữ liệu người dùng trái ý; xoá cache nếu cần.

## Chi phí vận hành & bảo trì

- **Tài nguyên**: CPU/GPU tùy mô hình; dùng biến thể “nhẹ” cho demo để tối ưu chi phí.
- **Phụ thuộc gói**: ghi rõ phiên bản trong **requirements.txt** để tránh lỗi build khi triển khai.
- **Giới hạn**: demo không thay thế hạ tầng production (monitoring, canary, rollback).

## Checklist ra quyết định “có dùng Gradio không”

Mục tiêu là **demo/POC**, không phải sản phẩm production dài hạn.

Có thể mô tả bài toán với **Input–Model–Output** rõ ràng.

UI chỉ cần **form đơn giản** (upload, slider, nút chạy) và vùng hiển thị kết quả.

Chấp nhận cơ chế **chia sẻ tạm** hoặc deploy nhẹ (Spaces).

## Kết nối sang phần kế tiếp

Sau khi nắm lý do sử dụng Gradio, phần tiếp theo sẽ trình bày **khung tư duy Input–Model–Output** chi tiết hơn và cách ánh xạ trực tiếp sang các component của Gradio để bạn bắt tay vào dựng UI đúng ngay từ đầu.



### 3 Khung tư duy: Input – Model – Output

#### Mục tiêu của phần này

Xây một demo vững chắc bắt đầu từ việc **định nghĩa hợp đồng dữ liệu** (data contract) giữa UI và hàm suy luận. Phần này giúp bạn:

- Xác định chính xác **Input** người dùng đưa vào.
- Đặt chuẩn **Model function** (chữ ký hàm, kiểu dữ liệu, xử lý lỗi).
- Thiết kế **Output** dễ đọc, dễ kiểm thử, ít bất ngờ.

#### 3.1 Đặt “hợp đồng dữ liệu” ngay từ đầu

Một hợp đồng dữ liệu nên chỉ rõ:

1. **Loại dữ liệu** (text, số, ảnh, tệp, danh sách...).
2. **Miền giá trị hợp lệ** (VD:  $\text{temperature} \in [0, 1]$ , kích thước ảnh  $\leq 2k \times 2k$ ).
3. **Ràng buộc** (bắt buộc/tuỳ chọn, số lượng tối đa tệp).
4. **Định dạng trả về** (chuỗi, JSON, ảnh đã vẽ, tệp ZIP...).

Lợi ích: giúp bạn chọn đúng component của Gradio, và viết mã suy luận ít “if-else” chứa chấy.

#### 3.2 Ánh xạ Input ↔ Component

Kiểu input	Component gợi ý	Ghi chú triển khai
Văn bản ngắn/dài	Textbox / TextArea	lines phù hợp; hỗ trợ .submit()
Số, tham số liên tục	Number, Slider	Dùng Slider kèm .release() cho tác vụ nặng
Lựa chọn rời rạc	Radio, Dropdown	Thiết lập choices, value mặc định
Ảnh	Image	type="pil" hoặc "numpy"; resize trước khi suy luận
Tệp bất kỳ	File	Kiểm tra kích thước/định dạng; dọn tệp tạm
Nhiều mục	CheckboxGroup, Gallery	Chú ý ép kiểu sang danh sách

#### 3.3 Thiết kế hàm model\_fn: chữ ký, kiểm tra, trả về

Mẫu chữ ký hàm khuyến nghị:

```

1 from typing import Tuple, Dict, Any, Optional
2 from PIL import Image
3 import numpy as np
4
5 def model_fn(
6     text: Optional[str] = None,
7     image: Optional[Image.Image] = None,
8     threshold: float = 0.5,
9     options: Optional[Dict[str, Any]] = None
10 ) -> Tuple[Any, Dict[str, Any]]:
11     """

```

```

12     Tra ve (primary_output, meta_info).
13     - primary_output: du lieu chinh de hien thi (str/Image/ndarray/JSON).
14     - meta_info: thong tin bo sung (thoi gian xu ly, tham so, canh bao...).
15     """
16     # 1) Validate
17     if (text is None) and (image is None):
18         return "Vui long nhap van ban hoac upload anh.", {"ok": False}
19
20     if not (0.0 <= threshold <= 1.0):
21         threshold = max(0.0, min(1.0, threshold))
22
23     # 2) Chuyen doi kieu (neu can)
24     if image is not None and not isinstance(image, Image.Image):
25         image = Image.fromarray(np.asarray(image))
26
27     # 3) Goi mo hinh / suy luan
28     # ... your inference here ...
29     result = {"message": "OK", "threshold": threshold}
30
31     # 4) Tao output chinh + meta
32     primary_output = str(result) # vi du: tra JSON dang chuoi
33     meta = {"ok": True, "params": {"threshold": threshold}}
34     return primary_output, meta

```

Listing 3: Chữ ký hàm suy luận và xử lý lỗi nhất quán

**Nguyên tắc:** luôn trả về cùng một *hình dạng dữ liệu* để UI hiển thị ổn định.

### 3.4 Tiền xử lý (preprocess) & Hậu xử lý (postprocess)

#### Tiền xử lý

- Chuẩn hoá text (strip, lower, bỏ ký tự lạ nếu cần).
- Với ảnh: convert sang RGB, resize tối đa, chuẩn hoá dtype.
- Kiểm tra kích thước tệp, số lượng tệp.

#### Hậu xử lý

- Định dạng điểm số, nhãn (f"score:.2f").
- Vẽ bbox/overlay lên ảnh (nếu là CV).
- Kèm **meta info**: thời gian xử lý, tham số đã dùng.

### 3.5 Thiết kế Output: dễ đọc, kiểm thử, và nhất quán

Một số mẫu output phổ biến:

1. **Chữ (textbox/markdown)**: ngắn gọn, có gạch đầu dòng, nêu điểm chính.
2. **Ảnh** đã vẽ sẵn: trả về PIL Image/ndarray đã annotate.
3. **Bảng/JSON**: dùng Dataframe/JSON để hiển thị cấu trúc.
4. **Tệp tải về**: zip/csv kết quả hàng loạt.

Quy tắc vàng:

- `interactive=False` cho phần chỉ để xem.
- Trả về **meta** (dictionary) để debug/test A/B.
- Giao diện **không thay đổi hình dạng** khi giá trị biên (VD: “không phát hiện gì”) — luôn hiển thị khung kết quả ổn định.

### 3.6 Ví dụ 1: NLP đơn giản (text + slider → text)

```

1 import gradio as gr
2 import time
3
4 def nlp_demo(text: str, temperature: float = 0.5):
5     start = time.time()
6     if not text or text.strip() == "":
7         return "Vui lòng nhập văn bản.", {"ok": False}
8     # Fake generation: đảo chuỗi làm ví dụ
9     out = f"[temp={temperature:.2f}] " + text[::-1]
10    meta = {"ok": True, "elapsed": round(time.time()-start, 3), "temp": temperature}
11    return out, meta
12
13 with gr.Blocks(title="NLP Demo") as demo:
14     with gr.Row():
15         with gr.Column(scale=1):
16             inp = gr.Textbox(label="Nhập văn bản")
17             temp = gr.Slider(0.0, 1.0, value=0.5, step=0.05, label="Temperature")
18             run = gr.Button("Chạy")
19         with gr.Column(scale=2):
20             out = gr.Textbox(label="Kết quả", interactive=False, lines=6)
21             meta = gr.JSON(label="Meta")
22     run.click(nlp_demo, inputs=[inp, temp], outputs=[out, meta])
23     temp.release(nlp_demo, inputs=[inp, temp], outputs=[out, meta]) % dùng .release
24     tránh gọi liên tục
25 demo.launch()

```

Listing 4: Input text, tham số slider; Output text + meta

### 3.7 Ví dụ 2: Ảnh (image + threshold → image + text)

```

1 from PIL import Image
2 import numpy as np
3 import gradio as gr
4
5 def fake_detect(image: Image.Image, thr: float = 0.3):
6     if image is None:
7         return None, "Vui lòng upload ảnh."
8     # Ví dụ: ‘vạch’ một strip màu lên ảnh để minh họa hậu xử lý
9     arr = np.array(image.convert("RGB"))
10    h, w, _ = arr.shape
11    arr[:max(5, int(h*thr)), :, :] = 255 # kẻ một vạch trắng theo threshold
12    annotated = Image.fromarray(arr)
13    info = f"Da annotate: strip={int(h*thr)}px (thr={thr:.2f})"
14    return annotated, info
15
16 with gr.Blocks(title="Image Demo") as app:
17     with gr.Row():
18         with gr.Column(scale=1):
19             img = gr.Image(label="Ảnh đầu vào", type="pil")
20             thr = gr.Slider(0.05, 0.95, value=0.3, step=0.05, label="Threshold")

```

```

21         run = gr.Button("Annotate")
22         with gr.Column(scale=2):
23             out_img = gr.Image(label="Anh da annotate", interactive=False)
24             out_txt = gr.Textbox(label="Thong tin", interactive=False)
25         run.click(fake_detect, inputs=[img, thr], outputs=[out_img, out_txt])
26         thr.release(fake_detect, inputs=[img, thr], outputs=[out_img, out_txt])
27 app.launch()

```

Listing 5: Input PIL Image; Output Image da annotate + thong tin tom tat

### 3.8 Nhiều Input/Output: quy ước thứ tự và kiểu

Khi có nhiều đầu vào/đầu ra, hãy đảm bảo:

- Thứ tự `inputs=[...]` và `outputs=[...]` khớp chữ ký hàm.
- Kiểu trả về đúng (VD: `Tuple[Image, str]`).
- Không thay đổi số lượng output theo trường hợp (kể cả khi lỗi).

```

1 def multi_io(text: str, image: Image.Image, k: int = 3):
2     if image is None or not text:
3         return None, "Thieu du lieu dau vao.", {"ok": False}
4     # ... suy luan ...
5     annotated = image # vi du
6     summary = f"Text={text[:20]}..., topK={k}"
7     meta = {"ok": True, "k": k}
8     return annotated, summary, meta
9
10 # outputs phai nhan 3 muc tuong ung
11 btn.click(multi_io, inputs=[inp_text, inp_img, inp_k], outputs=[out_img, out_txt,
    out_meta])

```

Listing 6: Nhieu input/output va tra ve on dinh

### 3.9 Xử lý lỗi và trạng thái biên (edge cases)

- Loại dữ liệu sai: hiển thị thông báo thay vì raise exception lan ra UI.
- Dữ liệu trống/quá lớn: trả thông điệp hướng dẫn (kích thước tối đa, định dạng hợp lệ).
- Không có kết quả: vẫn trả về khung output cố định (VD: ảnh gốc + “không phát hiện”).

### 3.10 Hiệu năng và ổn định

- Tiền xử lý nhẹ (resize ảnh), `.release` cho slider, `queue=True` khi nhiều người dùng.
- Tách phần nặng ra ngoài hàm sự kiện (load model một lần ở global scope).
- Cache kết quả trung gian nếu hợp lý (nhưng cân nhắc bộ nhớ).

### 3.11 Bảo mật dữ liệu (ở mức demo)

- Hạn chế lưu trữ dữ liệu người dùng; dọn tệp tạm sau khi xử lý.
- Không log thông tin nhạy cảm; nêu rõ phạm vi sử dụng ở phần mô tả.
- Khi chia sẻ liên kết, nhắc người dùng không đưa dữ liệu cá nhân.

## Tóm tắt

“Input – Model – Output” là xương sống của mọi demo Gradio. Khi hợp đồng dữ liệu rõ ràng, bạn sẽ:

- Chọn component chính xác, tránh đổi kiểu nhiều lần.
- Viết hàm suy luận sạch, dễ kiểm thử.
- Hiển thị kết quả nhất quán, thân thiện với người dùng.

Phần tiếp theo sẽ áp dụng khung này để so sánh và chọn giữa **Interface** và **Blocks**.

## 4 Các component hay dùng

### Mục tiêu của phần này

Chọn đúng component giúp:

- Giảm thời gian ghép nối giữa UI và hàm suy luận.
- Hạn chế lỗi ép kiểu (`PIL`  $\rightarrow$  `NumPy`, `str`  $\rightarrow$  `float` ...).
- Tạo trải nghiệm sử dụng tự nhiên với ít tùy chọn nhưng đủ.

### 4.1 Nguyên tắc chọn component

1. **Trung thực với kiểu dữ liệu:** nhập số dùng `Slider/Number`; ảnh dùng `Image`; tệp dùng `File`.
2. **Giới hạn sớm:** đặt ràng buộc tại UI (`min/max`, `step`, `choices`) để giảm nhánh kiểm tra ở backend.
3. **Giảm tự do nếu không cần:** dùng `Radio/Dropdown` thay vì `Textbox` tự do cho tham số phân loại.
4. **Output chỉ để xem:** thêm `interactive=false` để tránh người dùng chỉnh nhầm.

### 4.2 Nhóm Input cơ bản

#### Textbox / TextArea

- Dùng cho văn bản tự do (`prompt`, câu hỏi, mô tả).
- Thuộc tính gợi ý: `lines`, `max_lines`, `placeholder`.
- Gợi ý sự kiện: `.submit()` để nhấn Enter chạy luôn.

#### Number / Slider

- Dùng cho số/tham số liên tục (`confidence`, `temperature`).
- Thuộc tính: `minimum`, `maximum`, `step`, `value` (mặc định).
- Với tác vụ nặng: ưu tiên `.release()` thay vì `.change()` khi kéo `Slider`.

#### Radio / Dropdown / Checkbox / CheckboxGroup

- Thích hợp cho lựa chọn ràng buộc (một hoặc nhiều).
- Đặt `choices` (danh sách), `value` (mặc định).
- Tương ứng ở back-end: map sang `Enum` hoặc `set` để xử lý rõ ràng.

#### File

- Dùng khi loại tệp không cố định (`CSV`, `ZIP`, `PDF`...).
- Kiểm tra kích thước/đuôi tệp trong hàm suy luận; dọn tệp tạm nếu có ghi ra đĩa.

**Image**

- Chọn `type="pil"` để nhận `PIL.Image` hoặc `type="numpy"` để nhận `ndarray`.
- Nên **resize** ảnh lớn (ví dụ cạnh dài  $\leq 1280$ ) trước khi suy luận trên CPU.

**Audio / Video**

- Trả về đường dẫn tạm hoặc mảng mẫu tùy chế độ; cân nhắc thời lượng tối đa.
- Với Video, xem xét trích keyframe trước khi đưa vào mô hình nếu inference nặng.

**4.3 Nhóm Output phổ biến****Textbox / Label / Markdown**

- Dùng cho kết quả ngắn gọn, chú thích, hoặc định dạng nhẹ (Markdown).
- Đặt `interactive=false` cho output để tránh sửa nhầm.

**Image / Gallery**

- Cho thị giác máy tính (ảnh đã annotate, segmentation overlay).
- Gallery phù hợp hiển thị nhiều ảnh (top-k, trước/sau).

**JSON / Dataframe / HTML**

- JSON: debug thông số, trả cấu trúc kết quả.
- Dataframe: bảng kết quả, thống kê.
- HTML: mô tả giàu định dạng; thận trọng với nội dung người dùng nhập vào (escape).

**File (tải về)**

- Gộp kết quả hàng loạt thành ZIP/CSV cho người dùng tải.
- Đặt tên tệp rõ ràng, có timestamp để phân biệt phiên chạy.

**4.4 Mẫu ánh xạ nhanh Input  $\rightarrow$  Output**

Bài toán	Input	Output
Phân tích cảm xúc (NLP)	Textbox (câu)	Label + JSON (score)
Tóm tắt văn bản	Textbox (đoạn)	Textbox/Markdown (tóm tắt)
Object Detection	Image + Slider (conf)	Image (annotate) + Textbox (tóm tắt)
Phân lớp ảnh	Image	Label + JSON (top-k)
OCR	Image/File (PDF)	Textbox (văn bản) + File (.txt)
Batch inference	File (CSV/ZIP)	File (CSV/ZIP kết quả)

#### 4.5 Ví dụ 1: Textbox + Slider → Textbox + JSON

```

1 import gradio as gr, time, random
2
3 def nlp_infer(text: str, temperature: float = 0.5):
4     t0 = time.time()
5     if not text or text.strip()=="":
6         return "Vui long nhap van ban.", {"ok": False}
7     # minh hoa: dao chu va chen nhieu tu theo temperature
8     words = text.split()[::-1]
9     k = max(1, int(len(words)*temperature))
10    out = " ".join(words[:k])
11    meta = {"ok": True, "elapsed": round(time.time()-t0,3), "temp": temperature, "len": len(words)}
12    return out, meta
13
14 with gr.Blocks(title="NLP Example") as demo:
15     with gr.Row():
16         with gr.Column(scale=1):
17             inp = gr.Textbox(label="Nhap van ban", lines=4, placeholder="Vi du: Xin chao ...")
18             temp = gr.Slider(0.0, 1.0, value=0.5, step=0.05, label="Temperature")
19             run = gr.Button("Chay")
20             with gr.Column(scale=2):
21                 out = gr.Textbox(label="Ket qua", interactive=False, lines=6)
22                 meta = gr.JSON(label="Meta")
23             run.click(nlp_infer, inputs=[inp, temp], outputs=[out, meta])
24             temp.release(nlp_infer, inputs=[inp, temp], outputs=[out, meta])
25 demo.launch()

```

Listing 7: NLP: tham so temperature va ket qua dang text + JSON

#### 4.6 Ví dụ 2: Image + Slider → Image + Textbox

```

1 import gradio as gr
2 from PIL import Image
3 import numpy as np
4
5 def annotate(image: Image.Image, thr: float = 0.3):
6     if image is None:
7         return None, "Vui long upload anh."
8     arr = np.array(image.convert("RGB"))
9     h, w, _ = arr.shape
10    # ve 1 thanh ngang ti le theo thr -> minh hoa hau xu ly
11    hh = max(6, int(h*thr))
12    arr[:hh, :, :] = [255, 255, 255]
13    return Image.fromarray(arr), f"Da ve strip {hh}px (thr={thr:.2f})"
14
15 with gr.Blocks(title="Image Example") as app:
16     with gr.Row():
17         with gr.Column(scale=1):
18             img = gr.Image(label="Anh dau vao", type="pil")
19             thr = gr.Slider(0.05, 0.95, value=0.3, step=0.05, label="Threshold")
20             run = gr.Button("Annotate")
21             with gr.Column(scale=2):
22                 out_img = gr.Image(label="Anh da annotate", interactive=False)
23                 out_txt = gr.Textbox(label="Thong tin", interactive=False, lines=4)
24             run.click(annotate, inputs=[img, thr], outputs=[out_img, out_txt])
25             thr.release(annotate, inputs=[img, thr], outputs=[out_img, out_txt])

```



26 `app.launch()`

Listing 8: CV: annotate gia lap theo threshold

#### 4.7 Xử lý file và media an toàn

- **Giới hạn kích thước:** từ chối tệp quá lớn (thông báo rõ ràng).
- **Kiểm tra đuôi tệp:** chỉ nhận các định dạng bạn hỗ trợ (CSV/JPG/PNG ...).
- **Dọn tệp tạm:** nếu lưu ra ổ đĩa, xóa tệp sau khi trả kết quả.
- **Resize/convert sớm:** chuẩn hoá ảnh (RGB, kích thước) để giảm lỗi downstream.

#### 4.8 Mẹo UI/UX nhỏ nhưng hiệu quả

- **Placeholder** mô tả rõ *đầu vào hợp lệ*.
- **Giá trị mặc định** hợp lý: **Slider** bắt đầu ở giá trị thường dùng.
- **Nút Clear/Sample:** cho phép thử nhanh và reset trạng thái.
- **Trạng thái rỗng:** khi chưa có kết quả, vẫn hiển thị khung output để tránh “nhảy UI”.

#### Tóm tắt

Chọn component đúng giúp demo mạch lạc, ít lỗi ép kiểu và tối ưu chi phí tính toán. Hãy:

1. Bắt đầu từ kiểu dữ liệu **thật** mà mô hình cần.
2. Giới hạn sớm bằng **choices**, **min/max**, **step**.
3. Sử dụng **interactive=false** cho output và thêm **Clear/Sample**.

Phần kế tiếp sẽ so sánh cụ thể giữa **Interface** và **Blocks**, kèm tiêu chí chọn và ví dụ “đổi cỡ” từ bản đơn giản sang bản linh hoạt.

## 5 Hai cách dựng giao diện: Interface vs. Blocks

### Mục tiêu của phần này

Làm rõ khi nào dùng **Interface** để *lên demo thật nhanh*, khi nào dùng **Blocks** để *kiểm soát UI/UX và luồng tương tác phức tạp*. Đồng thời, đưa ra lộ trình **nâng cấp** (migrate) từ Interface → Blocks mà *không phải viết lại logic suy luận*.

### 5.1 Interface: đường tắt cho MVP/POC

#### Khi nào dùng?

- Bài toán có **một hàm suy luận trung tâm** (input → output).
- Tham số ít, thao tác tuyến tính: nhập → bấm nút → xem kết quả.
- Mục tiêu là **có link chạy được ngay** để trình diễn/thu phản hồi.

#### Đặc trưng

- Khai báo ngắn gọn: `fn`, `inputs`, `outputs`.
- Tự sinh UI với layout tối thiểu; ít phải quan tâm bố cục.
- Dễ đọc mã, ít “keo dán” giữa UI và backend.

```

1 import gradio as gr
2
3 def score_sentiment(text: str, threshold: float = 0.5):
4     if not text or text.strip()=="":
5         return "Vui lòng nhập văn bản."
6     # minh hoa: diem > threshold -> "pos" ngược lại "neg"
7     s = (sum(map(ord, text)) % 101) / 100.0
8     label = "positive" if s >= threshold else "negative"
9     return f"score={s:.2f} -> {label} (thr={threshold:.2f})"
10
11 demo = gr.Interface(
12     fn=score_sentiment,
13     inputs=[gr.Textbox(label="Văn bản"), gr.Slider(0,1,0.5,0.01,label="Threshold")],
14     outputs=gr.Textbox(label="Kết quả", interactive=False),
15     title="Sentiment MVP",
16     description="Nhập văn bản, điều chỉnh threshold, nhận kết quả."
17 )
18 demo.launch()
```

Listing 9: Interface: demo một-hàm

#### Ưu/nhược

- **Ưu:** Nhanh, ít mã, gần như không thể sai layout; rất hợp để dạy/hackathon/POC.
- **Nhược:** Khó làm **nhiều luồng** (multi-step), bố cục phức tạp, hoặc nhiều sự kiện tinh chỉnh.

## 5.2 Blocks: “lego” để mở rộng UI/UX và luồng

### Khi nào dùng?

- Cần **bố cục** (hàng/cột/tab/accordion), widget phụ trợ (Clear, Sample, Download).
- Có **nhiều sự kiện** (click, change, submit, release) và **nhiều hàm** xử lý.
- Muốn giữ **trạng thái** giữa các thao tác, hoặc tách nhiều pipeline.

### Đặc trưng

- Dựng UI theo ngữ nghĩa: `with gr.Blocks(): with gr.Row(): with gr.Column(): ...`
- **Kết nối sự kiện** chi tiết: `.click()`, `.change()`, `.submit()`, `.release()` ...
- Hỗ trợ tốt việc **tổ chức nhiều hàm suy luận** và phối hợp kết quả.

```

1 import gradio as gr
2
3 def preprocess(x):
4     return x.strip()
5
6 def infer(x, t):
7     if not x: return "Trong.", {"ok": False}
8     score = (sum(map(ord, x)) % 101) / 100.0
9     label = "positive" if score >= t else "negative"
10    return f"score={score:.2f} -> {label}", {"ok": True, "thr": t, "len": len(x)}
11
12 def clear_all():
13     return "", "", {}
14
15 with gr.Blocks(title="Sentiment Advanced", queue=True) as app:
16     gr.Markdown("## Sentiment (Blocks)")
17     with gr.Row():
18         with gr.Column(scale=1):
19             raw = gr.Textbox(label="Van ban", lines=4)
20             thr = gr.Slider(0,1,0.5,0.01,label="Threshold")
21             with gr.Row():
22                 btn_clean = gr.Button("Preprocess")
23                 btn_run = gr.Button("Infer")
24                 btn_clear = gr.Button("Clear")
25             with gr.Column(scale=2):
26                 cleaned = gr.Textbox(label="Sau preprocess", interactive=False)
27                 output = gr.Textbox(label="Ket qua", interactive=False, lines=4)
28                 meta = gr.JSON(label="Meta")
29
30     btn_clean.click(preprocess, inputs=raw, outputs=cleaned)
31     btn_run.click(infer, inputs=[cleaned, thr], outputs=[output, meta])
32     thr.release(infer, inputs=[cleaned, thr], outputs=[output, meta])
33     btn_clear.click(clear_all, outputs=[cleaned, output, meta])
34
35 app.launch()

```

Listing 10: Blocks: bo cuc + nhieu su kien

## Ưu/nhược

- **Ưu:** Rất linh hoạt; dễ tạo **trải nghiệm giàu tương tác** và quản lý nhiều bước.
- **Nhược:** Nhiều mã hơn Interface; cần suy nghĩ **kiến trúc sự kiện** để tránh rối.

## 5.3 Tiêu chí chọn nhanh

Tình huống	Nên chọn	Lý do
Một hàm suy luận, 1–2 tham số	Interface	Lên demo tức thì, ít mã, khó sai
Cần bố cục 2 cột + nút Clear/Sample	Blocks	Control layout + nhiều sự kiện
Nhiều pipeline/luồng (tiền xử lý, suy luận, hậu xử lý)	Blocks	Tách hàm, kết nối sự kiện minh bạch
Prototype cực nhanh để xin ý kiến	Interface (sau đó nâng cấp)	Ra link nhanh, rồi migrate dần

## 5.4 Mẫu nâng cấp: từ Interface → Blocks

```
1 demo = gr.Interface(fn=infer, inputs=[gr.Textbox(), gr.Slider(0,1,0.5,0.01)],
2                      outputs=gr.Textbox())
3 demo.launch()
```

```
1 def preprocess(x): return x.strip()
2 def infer(x, t): # giu nguyen chu ky
3     # ... suy luận ...
4     return "ket qua"
```

```
1 with gr.Blocks() as app:
2     with gr.Row():
3         with gr.Column(scale=1):
4             raw = gr.Textbox()
5             thr = gr.Slider(0,1,0.5,0.01)
6             clean = gr.Button("Preprocess")
7             run = gr.Button("Infer")
8         with gr.Column(scale=2):
9             x = gr.Textbox(label="Sau preprocess", interactive=False)
10            y = gr.Textbox(label="Ket qua", interactive=False)
11        clean.click(preprocess, raw, x)
12        run.click(infer, [x, thr], y)
13 app.launch()
```

**Điểm mấu chốt:** không thay đổi chữ ký infer; chỉ thay UI gọi vào đúng chỗ. Nhờ đó, việc nâng cấp không đụng tới logic mô hình.

## 5.5 Các mẫu thiết kế (patterns) thường gặp

### Bước 3: Blocks với bố cục + sự kiện

**Pattern 1: “Form trái – Kết quả phải”**

- Trái: input + tham số + nút hành động.
- Phải: kết quả (image/text) + meta (JSON/log).

**Pattern 2: “Tab theo chức năng”**

- Tab 1: Upload/nhập liệu.
- Tab 2: Kết quả/biểu đồ.
- Tab 3: Cấu hình nâng cao/Logs.

**Pattern 3: “Pipeline nhiều bước”**

- Nút Preprocess → Infer → Postprocess/Export.
- Mỗi bước một hàm; kết quả bước trước là input bước sau.

**5.6 Bẫy thường gặp và cách tránh**

- **Gọi model quá nhiều lần:** tránh `.change` với Slider cho tác vụ nặng, dùng `.release`.
- **Không giữ định dạng output:** luôn trả về cùng kiểu/shape để UI không “nhảy”.
- **Trộn logic vào UI:** tách hàm `preprocess/infer/postprocess` riêng, UI chỉ “điều phối”.
- **Thiếu nút Clear/Sample:** khó kiểm thử nhanh, khó reset trạng thái.

**5.7 Hiệu năng & khả năng mở rộng**

- Đặt tải mô hình (`load weights`) ở *global scope* để không nạp lại mỗi lần.
- Dùng `queue=True` khi chờ nhiều người dùng; cân nhắc batch nếu phù hợp.
- Với ảnh/video lớn: `resize`/giới hạn độ dài để giảm độ trễ.

**Tóm tắt**

Interface cho tốc độ và sự đơn giản; Blocks cho kiểm soát và mở rộng. Lộ trình hợp lý là:

1. **Bắt đầu** với Interface để có demo chạy ngay.
2. **Giữ** logic suy luận độc lập với UI.
3. **Nâng cấp** lên Blocks khi cần bố cục/sự kiện phức tạp, mà không đổi chữ ký hàm.

Ở phần tiếp theo, ta sẽ **bố cục UI** bằng Row/Column, Tab, Accordion để tạo trải nghiệm gọn gàng và dễ dùng.

## 6 Bố cục UI: Row/Column, Tab, Accordion

### Mục tiêu của phần này

Tổ chức giao diện sao cho:

- **Nhìn thấy ngay** các thao tác chính (nhập liệu, tham số, nút chạy).
- **Kết quả nổi bật**, không bị lẫn trong phần cấu hình.
- Có **không gian mở rộng** (nhiều chức năng, nhật ký, cấu hình nâng cao) mà không làm rối mắt.

### 6.1 Nguyên tắc bố cục tổng quát

1. **1 màn hình = 1 mục tiêu**: trang chính chỉ nên phục vụ một *luồng sử dụng* ưu tiên.
2. **Tách khu nhập liệu và kết quả**: bố trí “*Form bên trái – Kết quả bên phải*” cho hầu hết demo suy luận.
3. **Giảm chuyển động UI**: hạn chế thay đổi kích thước/thứ tự phần tử khi có kết quả mới; giữ khung kết quả cố định.
4. **Tối thiểu cuộn**: ưu tiên Row/Column thay vì xếp dọc quá dài; dùng Tab/Accordion cho phần phụ.

### 6.2 Row/Column: xương sống của layout

**Khái niệm** Row sắp xếp phần tử theo hàng ngang; Column sắp xếp theo cột dọc. Có thể lồng nhau để tạo bố cục hai cột, nhiều hàng.

**Tỉ lệ cột với scale** scale xác định tỉ lệ phân chia không gian trong một Row. Ví dụ: trái 1, phải 2 (1:2).

```

1 import gradio as gr
2
3 with gr.Blocks(title="Layout Co Ban") as demo:
4     gr.Markdown("## Demo: Form trai -- Ket qua phai")
5     with gr.Row():
6         # C t t r i (scale=1): n h p l i u v t h a m s
7         with gr.Column(scale=1):
8             inp = gr.Textbox(label="Nhap van ban", lines=4, placeholder="Vi du...")
9             thr = gr.Slider(0, 1, value=0.5, step=0.05, label="Threshold")
10            with gr.Row():
11                btn_run = gr.Button("Chay")
12                btn_clear = gr.Button("Clear")
13
14            # C t p h i (scale=2): k i t q u v m e t a
15            with gr.Column(scale=2):
16                out = gr.Textbox(label="Ket qua", interactive=False, lines=8)
17                meta = gr.JSON(label="Meta")
18
19            # G n s k i n m i n h h o
20            btn_clear.click(lambda: ("", {"note": "cleared"}), outputs=[out, meta])
21 demo.launch()
```

Listing 11: Pattern cơ bản: Form trái – Kết quả phải

**Chồng nhiều hàng trong một cột** Khi cần nhóm tham số, đặt chúng trong một Column riêng để tránh dồn thành “dải nút”.

```
1 with gr.Blocks(title="Rows trong Column") as app:
2     with gr.Row():
3         with gr.Column(scale=1):
4             with gr.Row():
5                 inp = gr.Textbox(label="Nhập 1", lines=3)
6             with gr.Row():
7                 thr = gr.Slider(0,1,0.5,0.05,label="Threshold")
8             with gr.Row():
9                 btn = gr.Button("Run")
10        with gr.Column(scale=2):
11            out = gr.Textbox(label="Kết quả", interactive=False, lines=6)
12        btn.click(lambda x, t: f"Input len={len(x)}; thr={t}", [inp, thr], out)
13 app.launch()
```

Listing 12: Nhiều hàng trong cột trái

### 6.3 Tab: phân tách chức năng theo ngữ nghĩa

Dùng Tab khi có nhiều *nhóm* chức năng: ví dụ Upload/Infer, Phân tích kết quả, Cấu hình nâng cao, Logs.

```
1 with gr.Blocks(title="Tab Layout") as demo:
2     with gr.Tab("Chức năng"):
3         with gr.Row():
4             with gr.Column(scale=1):
5                 img = gr.Image(label="Ảnh đầu vào", type="pil")
6                 conf = gr.Slider(0.05, 0.95, value=0.25, step=0.05, label="Confidence")
7                 btn = gr.Button("Detect")
8             with gr.Column(scale=2):
9                 out_img = gr.Image(label="Ảnh đã annotate", interactive=False)
10                out_txt = gr.Textbox(label="Thông tin", interactive=False, lines=8)
11            btn.click(lambda im, c: (im, f"conf={c:.2f}"), [img, conf], [out_img, out_txt])
12
13        with gr.Tab("Kết quả phân tích"):
14            gr.Markdown("### Biểu đồ/Thống kê (placeholder)")
15            stats = gr.JSON(label="Thống kê")
16
17        with gr.Tab("Cấu hình nâng cao"):
18            adv = gr.CheckboxGroup(choices=["A","B","C"], label="Tùy chọn")
19            save = gr.Button("Lưu cấu hình")
20            save.click(lambda x: {"saved": x}, adv, stats)
21
22        with gr.Tab("Logs"):
23            logs = gr.Textbox(label="Nhật ký", lines=12, interactive=False)
24 demo.launch()
```

Listing 13: Bố cục theo Tab: Chức năng, Kết quả, Cấu hình

### 6.4 Accordion: giấu bớt phần ít dùng

Accordion phù hợp với phần cấu hình hiếm khi dùng tới (ví dụ: seed, batch size, định dạng xuất).

```
1 with gr.Blocks(title="Accordion") as demo:
2     with gr.Row():
3         with gr.Column(scale=1):
4             text = gr.Textbox(label="Nhập văn bản", lines=3)
5             run = gr.Button("Chạy")
```

```

6         with gr.Accordion("Cau hinh nang cao", open=False):
7             seed = gr.Number(value=42, label="Seed")
8             top_k = gr.Slider(1, 20, value=5, step=1, label="Top-K")
9         with gr.Column(scale=2):
10             out = gr.Textbox(label="Ket qua", interactive=False, lines=8)
11
12     def infer(x, s, k):
13         return f"text='{x[:20]}...', seed={int(s)}, top_k={int(k)}"
14
15     run.click(infer, [text, seed, top_k], out)
16 demo.launch()

```

Listing 14: Accordion cho cau hinh hiem dung

## 6.5 Tổ chức “sidebar” tham số

Với nhiều tham số, có thể mô phỏng một “sidebar” bằng cột trái hẹp (scale nhỏ) chứa toàn bộ slider/checkbox.

```

1 with gr.Blocks(title="Sidebar") as demo:
2     with gr.Row():
3         with gr.Column(scale=1, min_width=280):
4             gr.Markdown("### Tham so")
5             temp = gr.Slider(0,1,0.5,0.01,label="Temperature")
6             top_p = gr.Slider(0,1,0.9,0.01,label="Top-p")
7             rep = gr.Slider(0,2,1.0,0.05,label="Repetition penalty")
8             run = gr.Button("Generate")
9         with gr.Column(scale=2):
10             prompt = gr.Textbox(label="Prompt", lines=6)
11             output = gr.Textbox(label="Output", interactive=False, lines=12)
12             run.click(lambda p, t, tp, r: f"[t={t:.2f}, tp={tp:.2f}, r={r:.2f}] {p}",
13                       [prompt, temp, top_p, rep], output)
14 demo.launch()

```

Listing 15: Sidebar tham so ben trai

## 6.6 Chia nhóm “kết quả chính” và “siêu dữ liệu”

Hiển thị kết quả trực quan (ảnh/nhân) ở vị trí nổi bật; phần *siêu dữ liệu* (JSON, thời gian xử lý, thông số) ở khu vực phụ để không làm nhiễu.

```

1 with gr.Blocks(title="Result + Meta") as app:
2     with gr.Row():
3         with gr.Column(scale=2):
4             main = gr.Image(label="Ket qua chinh", interactive=False)
5         with gr.Column(scale=1):
6             meta = gr.JSON(label="Thong tin bo sung")
7     # minh hoa
8     app.load(lambda: ("https://picsum.photos/800", {"elapsed": 0.12, "version": "demo"}),
9             inputs=None, outputs=[main, meta])
10 app.launch()

```

Listing 16: Ket qua chinh + sieu du lieu

## 6.7 Anti-patterns cần tránh

- **Một cột dọc rất dài:** người dùng phải cuộn liên tục để tìm nút/kết quả; thay vào đó chia 2 cột hoặc dùng Tab.
- **Kết quả chen giữa form:** làm người dùng mất định hướng; nên để kết quả ở cột/phần riêng.



- **UI nhảy kích thước:** khi không có kết quả thì trống hoàn toàn, có kết quả thì chiếm chỗ lớn; nên đặt container cố định (ví dụ `lines` cho Textbox, hoặc sẵn khung `Image`).
- **Quá nhiều tham số cấp 1:** nhét hết ra ngoài khiến form rối; gom nhóm vào Accordion/Tab.

## 6.8 Mẹo vi mô để UI gọn gàng

- Dùng `min_width` cho cột chứa nhiều slider để tránh co cụm.
- Giới hạn `lines` và `max_lines` cho Textbox để giữ chiều cao ổn định.
- Nhân (`label`) ngắn gọn, nhất quán kiểu chữ; thêm `placeholder` để hướng dẫn.
- Nhóm các nút liên quan trong cùng Row (VD: `Run/Clear/Sample`).

## Tóm tắt

Bố cục tốt giúp người dùng:

1. **Biết phải làm gì** (khu nhập liệu rõ ràng).
2. **Thấy được điều gì** (kết quả nổi bật, ổn định).
3. **Không bị nhiễu** (cấu hình/nhập ký ẩn trong Tab/Accordion).

Ở phần tiếp theo, ta sẽ đi sâu vào **Event listeners** (`.click`, `.change`, `.submit`, `.release`) và cách chọn đúng để app mượt, tránh gọi suy luận thừa.

## 7 Event listeners: chọn đúng để app mượt

### Mục tiêu của phần này

Chọn đúng *sự kiện* để:

- Tránh gọi suy luận **thừa** (giảm độ trễ, tiết kiệm tài nguyên).
- Tăng **độ phản hồi** và tính dự đoán được của UI.
- Giữ **hình dạng output ổn định** trong mọi tình huống.

### 7.1 Bốn sự kiện cốt lõi

`.click()` Kích hoạt khi người dùng bấm nút. Phù hợp nhất cho “chạy suy luận” chủ động.

`.change()` Kích hoạt *mỗi khi* giá trị component thay đổi (nhập ký tự, kéo slider). Dễ gây *lặp gọi* nếu tác vụ nặng.

`.submit()` Kích hoạt khi nhấn Enter trong Textbox/TextArea. Trải nghiệm tự nhiên cho tác vụ text.

`.release()` Kích hoạt khi **thả chuột** khỏi Slider. Lý tưởng cho tham số ảnh hưởng mạnh đến chi phí tính toán (confidence, temperature...).

### 7.2 Quy tắc chọn nhanh

Tình huống	Sự kiện khuyến nghị
Nút “Chạy”, “Detect”, “Generate”	<code>.click()</code>
Nhập văn bản xong và nhấn Enter	<code>.submit()</code>
Điều chỉnh Slider cho tác vụ nặng (OD, sinh ảnh, LLM)	<code>.release()</code> (tránh <code>.change()</code> )
Những thay đổi nhẹ, chi phí thấp (lọc danh sách, hiển thị tức thì)	<code>.change()</code>

### 7.3 Mẫu kết nối sự kiện chuẩn

```

1 import gradio as gr
2
3 def run_infer(x, t):
4     # ... gọi mô hình ...
5     return f"done: len={len(x)}, thr={t:.2f}"
6
7 with gr.Blocks(title="Events Basic") as app:
8     txt = gr.Textbox(label="Nhập văn bản", lines=3)
9     thr = gr.Slider(0,1,0.5,0.01,label="Threshold")
10    btn = gr.Button("Chạy")
11    out = gr.Textbox(label="Kết quả", interactive=False, lines=4)
12
13    # Chạy bằng nút
14    btn.click(run_infer, inputs=[txt, thr], outputs=out)
15    # Chạy bằng Enter
16    txt.submit(run_infer, inputs=[txt, thr], outputs=out)
17    # Chạy khi thả slider
18    thr.release(run_infer, inputs=[txt, thr], outputs=out)

```

```
19
20 app.launch()
```

Listing 17: Kết nối sự kiện cơ bản cho một demo

## 7.4 Tránh lặp gọi với Slider

**Vấn đề** `.change()` trên Slider sẽ gọi hàm *liên tục* khi người dùng kéo — tệ cho tác vụ nặng (Object Detection, diffusion, LLM).

**Giải pháp** Dùng `.release()` để gọi *một lần* sau khi người dùng chọn xong giá trị.

```
1 conf = gr.Slider(0.05, 0.95, value=0.25, step=0.05, label="Confidence")
2 out_img = gr.Image()
3 out_txt = gr.Textbox()
4
5 def infer(image, conf):
6     # ... object detection ...
7     return image, f"conf={conf:.2f}"
8
9 conf.release(infer, inputs=[img, conf], outputs=[out_img, out_txt])
```

Listing 18: Ưu tiên `.release` cho tác vụ nặng

## 7.5 Ghép nhiều sự kiện cho một luồng

Có thể gán **nhiều sự kiện** vào cùng một hàm để hỗ trợ nhiều cách kích hoạt:

```
1 def pipeline(text, temp):
2     # ... inference ...
3     return f"[temp={temp:.2f}] {text[:-1]}"
4
5 btn.click(pipeline, [txt, temp], out)
6 txt.submit(pipeline, [txt, temp], out)
7 temp.release(pipeline, [txt, temp], out)
```

Listing 19: `.click` + `.submit` + `.release` tối cùng một hàm

## 7.6 Tách *preprocess/infer/postprocess* bằng sự kiện

```
1 def preprocess(x):
2     return x.strip()
3
4 def infer(x, t):
5     if not x: return "Trong."
6     # ... heavy model ...
7     return f"OK({t:.2f}): {x[:20]}..."
8
9 def postprocess(y):
10    return y.upper()
11
12 with gr.Blocks() as app:
13    raw = gr.Textbox(label="Raw", lines=3)
14    temp = gr.Slider(0,1,0.5,0.05,label="Temp")
15    btnP = gr.Button("Preprocess")
16    btnI = gr.Button("Infer")
17    btnO = gr.Button("Postprocess")
18
19    cleaned = gr.Textbox(label="Cleaned", interactive=False)
20    out1 = gr.Textbox(label="Infer out", interactive=False)
```

```

21 out2 = gr.Textbox(label="Post out", interactive=False)
22
23 btnP.click(preprocess, raw, cleaned)
24 btnI.click(infer, [cleaned, temp], out1)
25 btnO.click(postprocess, out1, out2)
26 app.launch()

```

Listing 20: Thiết kế pipeline nhiều bước qua sự kiện

## 7.7 Xử lý lỗi & trạng thái biên qua sự kiện

- **Luôn trả về cùng cấu trúc:** dù lỗi hay thành công ((output, meta)), UI sẽ không “nhảy”.
- **Hiển thị thông điệp thân thiện:** thay vì raise exception ra UI, trả chuỗi hướng dẫn (“Vui lòng upload ảnh.”).
- **Nút Clear:** gắn sự kiện trả về giá trị rỗng cho tất cả output liên quan.

```

1 def clear_all():
2     return "", {}, {} # ví dụ: out_text, out_log, meta
3
4 btn_clear.click(clear_all, outputs=[out_text, out_log, out_meta])

```

Listing 21: Nút Clear trả về trạng thái rỗng ổn định

## 7.8 Hàng đợi (queue=True) và phản hồi người dùng

- Bật queue=True khi dự kiến nhiều yêu cầu song song để điều phối lượt chạy.
- Cân nhắc hiển thị “đang xử lý” (progress/message) cho tác vụ dài (debug=True để thấy traceback khi lỗi).

```

1 with gr.Blocks(title="Queued App", queue=True) as demo:
2     # ... UI ...
3     pass
4 demo.launch(debug=True)

```

Listing 22: Bật hàng đợi cho app nhiều người dùng

## 7.9 Anti-patterns thường gặp

- **Gắn mọi thứ vào .change():** kéo/thả/nhập đều kích hoạt suy luận liên tục, gây lag.
- **Không ràng buộc tham số:** slider không có step/min/max dẫn tới giá trị “lạ”.
- **Output lúc có lúc không:** trả về số lượng output khác nhau giữa các nhánh (làm vỡ layout).

## Tóm tắt

1. Dùng .click() cho hành động có chủ đích, .submit() cho text, .release() cho slider tác vụ nặng.
2. Tránh .change() khi inference nặng; chỉ dùng cho thao tác nhẹ, tức thì.
3. Giữ cấu trúc trả về nhất quán và có **Clear** để reset nhanh.

Phần tiếp theo sẽ áp dụng toàn bộ vào một demo **Object Detection (YOLOv8n)** hoàn chỉnh, từ nạp model, viết hàm suy luận, dàn layout đến nối sự kiện.

## 8 Thực hành: Object Detection (YOLOv8n) với Gradio

### Mục tiêu của phần này

Xây một ứng dụng demo **Object Detection** hoàn chỉnh:

- Người dùng *upload ảnh* hoặc *tải ảnh mẫu*.
- Điều chỉnh **confidence** (Slider) và chạy suy luận.
- Nhận **ảnh đã annotate** (vẽ bounding box + nhãn) và **tóm tắt phát hiện**.

Ta sẽ dùng **YOLOv8n** (bản nano) từ thư viện **ultralytics** vì nhẹ, phù hợp demo trên CPU.

### 8.1 Chuẩn bị môi trường

Cài đặt gói cần thiết (Colab hoặc máy cá nhân):

```
1 pip install --upgrade gradio ultralytics opencv-python pillow torch torchvision
```

Listing 23: Cài đặt thư viện

### Ghi chú

- Nếu GPU có sẵn (CUDA), **torch** sẽ tự nhận; nếu không, demo vẫn chạy trên CPU với ảnh vừa phải.
- Với mạng chậm, lần đầu tải trọng số **yolov8n.pt** có thể mất thời gian.

### 8.2 Cấu trúc dự án gợi ý

```
.
app.py           # mã Gradio chính
requirements.txt # gói bắt buộc khi deploy (HF Spaces)
assets/          # (tùy chọn) ảnh minh họa, icon...
```

### 8.3 Hằng số và tiện ích xử lý ảnh

Ta viết một số tiện ích nhỏ để mã gọn, ổn định trên CPU:

```
1 import io, time, requests
2 from typing import Tuple, List, Optional, Dict, Any
3
4 import numpy as np
5 from PIL import Image
6
7 MAX_SIDE = 1280 # g i i h n k c h t h c n h g i m t r t r n
8                 CPU
9
10 SAMPLE_URL = "https://ultralytics.com/images/bus.jpg"
11
12 def load_sample_image() -> Image.Image:
13     resp = requests.get(SAMPLE_URL, timeout=10)
14     im = Image.open(io.BytesIO(resp.content)).convert("RGB")
15     return im
16
17 def ensure_rgb(im: Image.Image) -> Image.Image:
18     return im.convert("RGB") if im.mode != "RGB" else im
19
20 def resize_long_side(im: Image.Image, max_side: int = MAX_SIDE) -> Image.Image:
21     w, h = im.size
```

```

20 m = max(w, h)
21 if m <= max_side:
22     return im
23 scale = max_side / float(m)
24 new_w, new_h = int(round(w * scale)), int(round(h * scale))
25 return im.resize((new_w, new_h), Image.Resampling.LANCZOS)

```

Listing 24: Hằng số và tiện ích

## 8.4 Nạp mô hình một lần ở phạm vi toàn cục

Để tránh nạp lại mỗi lần người dùng bấm nút, ta đặt việc nạp model ở *global scope*.

```

1 from ultralytics import YOLO
2
3 # T i   m o d e l   ( l n   u   s í   t i   t r n g   s   v   c a c h e   n g   í   d   n g )
4 MODEL_NAME = "yolov8n.pt"
5 model = YOLO(MODEL_NAME)
6
7 # T y   c h n :   k i m   t r a   t h i   t   b
8 try:
9     import torch
10    DEVICE = "cuda" if torch.cuda.is_available() else "cpu"
11 except:
12    DEVICE = "cpu"

```

Listing 25: Nạp YOLOv8n một lần

## 8.5 Hàm suy luận: annotate ảnh + tóm tắt kết quả

Nguyên tắc:

- **Validate** đầu vào (ảnh trống).
- Chuẩn hoá: RGB, resize để giới hạn kích thước.
- Gọi `model.predict(..., conf=...)` rồi `plot()` để lấy ảnh annotate.
- Tạo bản tóm tắt từ boxes (class\_id, score, toạ độ).
- Trả về (annotated\_image, info\_text, meta) ổn định.

```

1 def detect_objects(image: Optional[Image.Image], conf: float = 0.25):
2     t0 = time.time()
3     if image is None:
4         return None, "Vui l ñ ng upload nh .", {"ok": False}
5
6     # T i n   x   l í
7     image = ensure_rgb(image)
8     image = resize_long_side(image, MAX_SIDE)
9     arr = np.array(image) # HWC RGB
10
11    # S u y   l u n
12    results = model.predict(arr, conf=conf, verbose=False, device=DEVICE)
13    res = results[0]
14
15    # V í   a n n o t a t e
16    annotated = res.plot() # ndarray RGB
17    annotated_pil = Image.fromarray(annotated)
18

```

```

19 # T m t t k i t q u
20 lines: List[str] = []
21 n = 0
22 if getattr(res, "boxes", None) is not None and len(res.boxes) > 0:
23     for b in res.boxes:
24         cls_id = int(b.cls[0])
25         score = float(b.conf[0])
26         x1, y1, x2, y2 = [float(x) for x in b.xyxy[0].tolist()]
27         lines.append(
28             f"- id={cls_id}, score={score:.2f}, "
29             f"box=({x1:.1f}, {y1:.1f}, {x2:.1f}, {y2:.1f})"
30         )
31     n = len(res.boxes)
32     info = "Ph t h i n :\n" + "\n".join(lines)
33 else:
34     info = "Kh ng ph t h i n i t ng n o ."
35
36 meta = {
37     "ok": True,
38     "elapsed_sec": round(time.time() - t0, 3),
39     "num_dets": n,
40     "conf": conf,
41     "device": DEVICE,
42     "input_size": list(arr.shape) # [H, W, C]
43 }
44 return annotated_pil, info, meta

```

Listing 26: Hàm suy luận Object Detection

## 8.6 Hàm tiện ích: Clear và tải ảnh mẫu

```

1 def clear_all():
2     # Tr v t r ng th i r ng cho c c output
3     return None, "", {}
4
5 def load_sample():
6     return load_sample_image()

```

Listing 27: Clear và Load Sample

## 8.7 Dựng UI bằng Blocks

Bố cục “Form trái – Kết quả phải”, có Sample/Detect/Clear, và confidence dùng `.release()` để tránh gọi liên tục.

```

1 import gradio as gr
2
3 with gr.Blocks(title="YOLOv8n Object Detection", queue=True) as demo:
4     gr.Markdown("## Object Detection Demo (YOLOv8n)")
5
6     with gr.Row():
7         # C t t r i: input + tham s + n t
8         with gr.Column(scale=1):
9             inp_image = gr.Image(label="nh u v o", type="pil")
10            conf = gr.Slider(0.05, 0.95, value=0.25, step=0.05, label="Confidence")
11            with gr.Row():
12                btn_sample = gr.Button("Load Sample")
13                btn_detect = gr.Button("Detect")
14                btn_clear = gr.Button("Clear")
15

```

```

16     # C t p h i : k i t q u + m e t a
17     with gr.Column(scale=2):
18         out_image = gr.Image(label=" n h a n n o t a t e", interactive=False)
19         out_info = gr.Textbox(label="Th n g t i n p h t h i n", interactive=False,
20                               lines=10)
21         out_meta = gr.JSON(label="Meta")
22
23     # S k i n
24     btn_sample.click(load_sample, outputs=inp_image)
25     btn_detect.click(detect_objects, inputs=[inp_image, conf], outputs=[out_image,
26                               out_info, out_meta])
27     btn_clear.click(clear_all, outputs=[out_image, out_info, out_meta])
28
29     # Tr n h g i s u y l u n l i n t c k h i k o s l i d e r :
30     conf.release(detect_objects, inputs=[inp_image, conf], outputs=[out_image,
31                               out_info, out_meta])
32
33 if __name__ == "__main__":
34     demo.launch(share=True, debug=True)

```

Listing 28: Gradio Blocks cho Object Detection

## 8.8 Kịch bản kiểm thử nhanh

- **Không upload ảnh:** bấm Detect → thông báo “Vui lòng upload ảnh.” (không lỗi).
- **Ảnh lớn (4K):** UI vẫn phản hồi, thời gian dài hơn; ảnh được `resize` về `MAX_SIDE`.
- **Confidence thấp/cao:** so sánh số lượng box và chất lượng phát hiện.
- **Nhấn Clear:** mọi output về trạng thái rỗng, layout ổn định.

## 8.9 Tối ưu hiệu năng

- **Resize ảnh:** giới hạn cạnh dài (`MAX_SIDE`) để giảm độ trễ, nhất là trên CPU.
- **Load model một lần:** đặt ở *global scope*.
- **Giảm số lần gọi:** dùng `.release` thay vì `.change` với Slider.
- **Chọn bản nhỏ:** `yolo8n.pt` thay vì `m/l/x` khi demo.

## 8.10 Xử lý lỗi và tình huống biên

- **Ảnh không hợp lệ:** bao try/except khi `open()` hoặc `np.array()`, trả thông điệp thân thiện.
- **Không phát hiện:** vẫn trả `annotated_image` (ảnh gốc) và chuỗi “Không phát hiện...” để UI không nhảy.
- **Thiếu internet (ảnh mẫu):** nếu tải `SAMPLE_URL` lỗi, thông báo “Không tải được ảnh mẫu.”.

## 8.11 An toàn dữ liệu (mức demo)

- Không lưu ảnh người dùng trừ khi thật cần; nếu có, nêu rõ mục đích và xoá sau khi xử lý.
- Không ghi log chi tiết chứa dữ liệu nhạy cảm.



### 8.12 Checklist hoàn tất phần thực hành

Chạy được trên CPU/GPU với ảnh mẫu và ảnh người dùng.

Output **ổn định**: (`image`, `info`, `meta`) trong mọi nhánh.

Nút `Sample/Detect/Clear` hoạt động như kỳ vọng.

Slider dùng `.release` để tránh suy luận thừa.

Ảnh lớn được `resize` trước khi suy luận.

### Kết nối sang phần tiếp theo

Ở phần kế tiếp, ta sẽ **đưa ứng dụng lên Hugging Face Spaces** để có *URL công khai bền vững*, kèm `requirements.txt` tối giản và mẹo tránh lỗi build.

## 9 Triển khai lên Hugging Face Spaces

### Mục tiêu của phần này

Đưa ứng dụng Gradio (đã chạy ổn cục bộ/Colab) lên **Hugging Face Spaces** để nhận *URL công khai bền vững*, tiện chia sẻ cho giảng viên, đồng đội, hoặc khách hàng.

### 9.1 Kiến thức nền tảng nhanh

- **Space** là một repo đặc biệt trên Hugging Face chứa mã nguồn ứng dụng. Bạn có thể chọn runtime Gradio.
- App tự build và chạy mỗi khi bạn **commit** (push) thay đổi.
- Có thể cấu hình **Secrets** (biến môi trường ẩn) và **Hardware** (CPU/GPU).

### 9.2 Chuẩn bị dự án

Cấu trúc tối thiểu:

```

.
├── app.py           # mã Gradio chính (điểm vào)
├── requirements.txt # phiên bản gói cần cài khi build
└── assets/         # (tùy chọn) ảnh minh họa, icon...
```

```

1 import gradio as gr
2
3 def echo(x: str) -> str:
4     return f"Echo: {x}"
5
6 demo = gr.Interface(
7     fn=echo,
8     inputs=gr.Textbox(label="Nhập gì đó"),
9     outputs=gr.Textbox(label="Két qua", interactive=False),
10    title="Hello Gradio on Spaces",
11    description="Demo tối thiểu"
12 )
13
14 if __name__ == "__main__":
15     demo.launch()
```

Listing 29: app.py tối giản cho Spaces

```

1 gradio>=4.0.0
2 ultralytics>=8.0.0
3 opencv-python>=4.8.0
4 pillow>=10.0.0
5 torch           # de Spaces tu chon ban phu hop voi phan c ng
6 torchvision
7 requests
```

Listing 30: requirements.txt tham khảo

**Mẹo:** khi đã có bản chạy ổn định, cần nhắc *ghim chặt* phiên bản (ví dụ `gradio==4.44.0`) để tránh build hỏng do thay đổi upstream.

### 9.3 Tạo Space và đẩy mã

#### requirements.txt (ghim phiên bản để build ổn định)

1. Đăng nhập Hugging Face → **Create new Space** → SDK: **Gradio**.
2. Chọn **Public** (mặc định) hoặc **Private** (nếu cần).
3. Tạo repo, sử dụng `git lfs` nếu có file lớn.
4. Thêm `app.py`, `requirements.txt`, và (tuỳ chọn) thư mục `assets/`. Commit và push.
5. Chờ tiến trình **Build** hoàn tất (~ vài phút). Khi thành công, Space sẽ hiển thị app và cấp URL công khai.

### 9.4 Cấu hình nâng cao

#### Hardware

- Mặc định: CPU. Với mô hình nặng (diffusion/LLM), chọn GPU phù hợp (tốn credit).
- Với OD dùng YOLOv8n, CPU thường đủ cho demo.

#### Secrets (biến môi trường)

- Dùng cho khóa API (ví dụ OpenAI, dịch vụ nội bộ).
- Truy cập ở phần **Settings** → **Repository secrets**; đọc trong mã qua `os.environ`.

#### Pin revision & cache

- Hugging Face có **cache** cho model/checkpoint (ví dụ `yolov8n.pt`).
- Có thể warm-up lần đầu lâu hơn; những lần sau nhanh hơn.

#### Startup script

- Với Gradio, không cần quy định port/tham số server; Spaces tự nhận.
- `demo.launch()` trong `app.py` là đủ.

### 9.5 Khắc phục sự cố thường gặp

- **Build fail do gói**: kiểm tra `requirements.txt`, ghim phiên bản cụ thể, loại bỏ gói không cần thiết.
- **Thiếu quyền mạng (tải mẫu)**: dùng `try/except` và thông báo “Không tải được ảnh mẫu.”; chuẩn bị ảnh fallback nội bộ.
- **Timeout khi nạp model lớn**: dùng bản nhẹ hơn, hoặc chuyển sang GPU nếu bắt buộc.
- **Vỡ UI sau cập nhật gói**: khôi phục phiên bản trước (ghi rõ version trong `requirements.txt`).

### 9.6 Checklist triển khai thành công

`app.py` chạy cục bộ được trước khi đẩy lên.

`requirements.txt` tối giản, có phiên bản phù hợp.

App hiện **online**, URL public truy cập ổn định.

Thêm **README** gắn ở Space mô tả cách dùng, dữ liệu mẫu, lưu ý bảo mật.

## Gợi ý mở rộng

- Thêm ảnh/video/README mô tả kết quả đầu ra.
- Sử dụng Tabs/Accordion để chia khu vực “Hướng dẫn”, “Giới thiệu mô hình”, “FAQ”.
- Bật `queue=True` cho app nhiều người dùng; thêm giới hạn vào/ra nếu cần.

## 10 Kết luận

### Tóm lược

Bằng cách áp dụng khung **Input – Model – Output** và lựa chọn đúng giữa **Interface / Blocks**, bạn có thể đưa một hàm suy luận hoặc mô hình ML thành **ứng dụng web** trong thời gian rất ngắn. Bố cục hợp lý (Row/Column, Tab, Accordion) và **sự kiện** phù hợp (`.click`, `.submit`, `.release`) giúp giảm độ trễ, ổn định UI, và nâng cao trải nghiệm người dùng.

### Giá trị thực tiễn

- **Demo nhanh/MVP**: chia sẻ đường link để lấy phản hồi sớm.
- **Giảng dạy/workshop**: giảm rào cản hạ tầng, tập trung vào logic học máy.
- **Nội bộ/POC**: kiểm chứng ý tưởng trước khi đầu tư xây dựng hệ thống production.

### Hướng đi tiếp theo

- Tách rõ các hàm `preprocess/infer/postprocess`, chuẩn hoá đầu vào/ra để dễ bảo trì.
- Thử nghiệm thêm **Blocks nâng cao**: *Tabs cho đa chức năng, Accordion cho cấu hình hiếm dùng.*
- Nâng cấp lên mô hình lớn hơn (khi cần) và cân nhắc **GPU** trên Spaces.
- Viết **README** rõ ràng + thêm hình minh hoạ kết quả để người dùng hiểu ngay cách sử dụng.

### Lời khuyên cuối

1. Bắt đầu nhỏ với **Interface**, có đường link **chạy được**.
2. Giữ logic mô hình *độc lập* với UI để **nâng cấp dần** lên **Blocks** mà không động đến lõi suy luận.
3. Luôn có **Clear**, **Sample**, và **thông báo lỗi thân thiện**.

Chúc bạn dựng demo thuận lợi và nhận được phản hồi hữu ích!