

Golang ile Generative AI Pratikleri

Gokonf '24

17 Şubat 2024
Yağmur YILDIZ

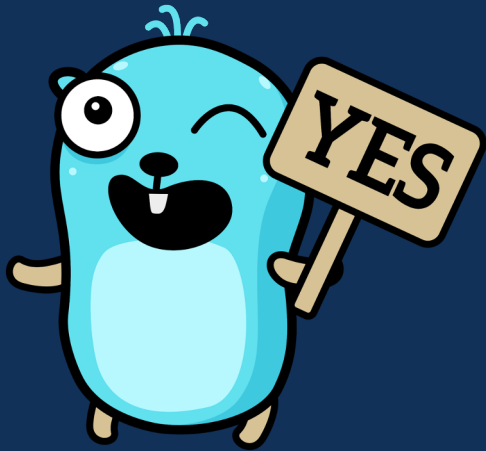


Akış

01 Neden Go & GenAI?

02 GenAI 101

03 Use Case: Txt2Img



Akış

01 Neden Go & GenAI?

02 GenAI 101

03 Use Case: Txt2Img

ÇIKARIM (INFERENCE)

- **Eğitimli yapay zeka modelini canlı canlı kullanarak çıktı elde etmek**
- **Bir model 30-40GB**
- **Bir cycle 5-20 sn**
- **Aynı anda birden fazla istek**



GO'NUN AVANTAJLARI

- **Performance**
 - **Efficiency**
 - **Concurrency**
 - **Scalability**
- + **Bulut tabanlı çözümlere için mükemmel olması**





Akış

01 Neden Go & GenAI?

02 GenAI 101

03 Use Case: Txt2Img

POPÜLER MODELLER

Kapalı Kaynak



ElevenLabs

Gemini



ANTHROPIC

Açık Kaynak



MISTRAL
AI_

BLM



LLaMA
by Meta



Cody

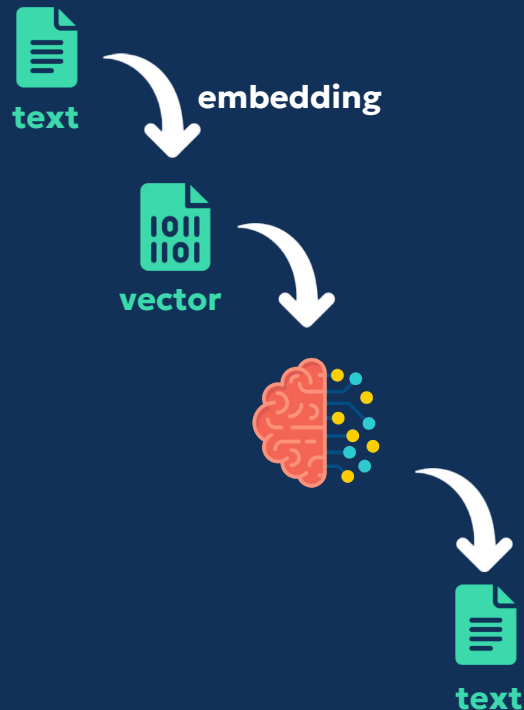
Meta AI
AudioCraft



BÜYÜK DİL MODELLERİ (LLMs)

Milyarlarca kelimelik büyük veriyle eğitilen derin öğrenme modelleridir.

- Transformers
- BERT (Bidirectional Encoder Representations Tr.)
- GPT (Gen Pre-Training Tr.)
- LaMDA (LM for Dialogue App.)



GÖRSEL ÜRETME

Daha çok diffusion modelleri kullanılan bu alanda amaç inputun öğelerini anlamlandırarak bir gürültünün öğeleri barındıran bir kompozisyona dönüşmesidir.

- Text2Image
- Image2Image
- Audio2Image
- Video2Image



Latent Noise



Generated Image

STABLE DIFFUSION MODELLERİ

Version	Release Date	Resolution	Parameters	Prompts Technology	Strengths	Weaknesses
1.4	08.2022	512x512	860M	CLIP	beginner-friendly, a little more artistic driven	long prompts, lower resolution
1.5	10. 2022	512x512	860M	CLIP	beginner-friendly, stronger portrait generation	long prompts, lower resolution
2.0	11.2022	768x768	-	OpenCLIP	shorter prompts, richer colors	aggressive NSFW filtering
2.1	12.2022	768x768	-	OpenCLIP	shorter prompts, richer colors	more “censored”, celeb filtered
XL 1.0	01.2023	1024x1024	3.5B	OpenCLIP & CLIP	shorter prompts, high resolution	requires GPU

SAMPLER, SCHEDULER, SEED

diversity, quality, speed, convergence

Sampler

Olasılık uzayında
modelin nasıl
çalışacağını belirler.
Euler, DDIM, DDPM

Scheduler

Modelin her sampledada
nasıl yakınsayacağını
belirler.
Linear, PNDM, Karras

Seed

Üretim sürecine
başlarken kullanılan
random değerdir.

SAMPLING STEP, CFG SCALE, SIZE

quality, creativity, computing power

Sampling Step

Samplerin kaç kez çalışacağını belirler. 30-75 arası tercih edilir.

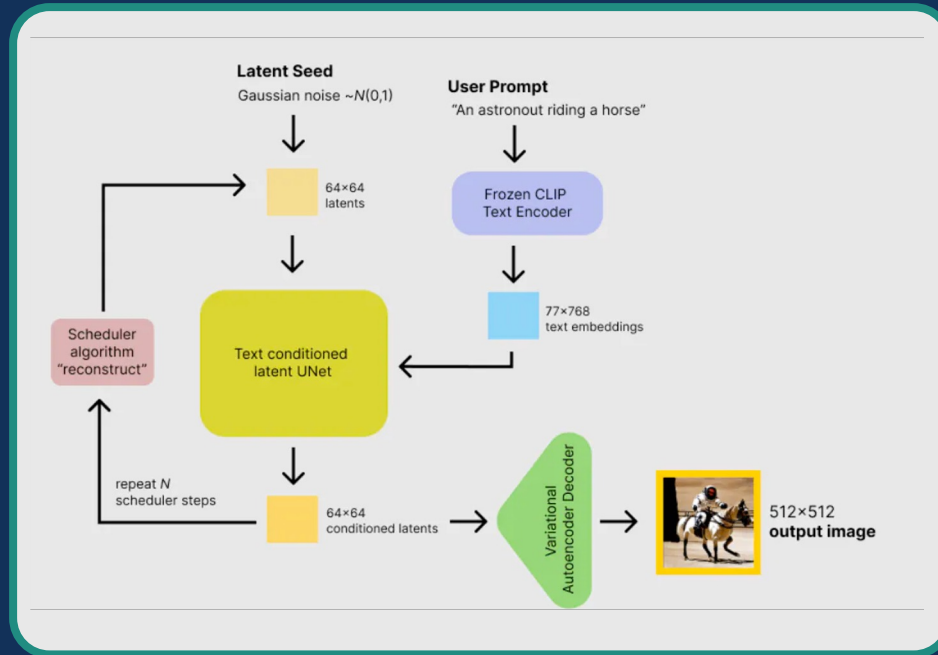
CFG Scale

Promptun generation üzerindeki etkisini belirler. 7-10 arası tercih edilir.

Size

Üretilcek resmin WxH boyutudur. Her modelin baz bir boyutu vardır.

STABLE DIFFUSION ALGORITHM





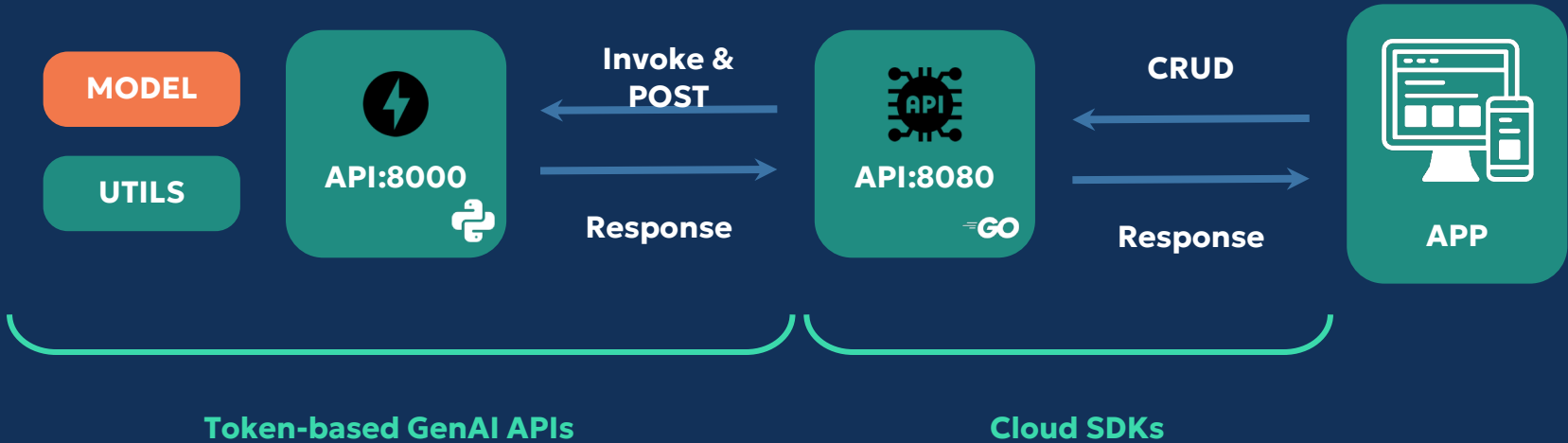
Akış

01 Neden Go & GenAI?

02 GenAI 101

03 Use Case: Txt2Img

TEMEL MİMARİ



GENERAL IMAGE REQUEST STRUCT

JSON

Bir prompt ile istediğimiz ölçülerde çıktı alabileceğimiz basit bir request yapısı

Constants

- Model
- Sampler
- Scheduler
- Output Format

```
1 {  
2   prompt: str,  
3   negative_prompt: str,  
4   num_inference_steps: int,  
5   width: int,  
6   height: int,  
7   guidance_scale: float  
8 }
```


MODEL API



Modeli barındıran repo ile konuşan ve modele requestlerin doğru atılmasını sağlayan katmandır.

- Python ile yapılması gereken işler
- Device settings
- Sampler, Scheduler, Seed

```
1
2 app = FastAPI()
3
4 @app.post("/txt2img/")
5 async def generate_image(request: ImageRequest):
6     try:
7         device = get_device_info()["device"]
8         dtype = get_device_info()["dtype"]
9         pipeline = create_pipeline("Realistic_Vision_V6.0_B1_noVAE", device, dtype)
10
11         generator = torch.Generator(device)
12         generator = torch.Generator(device=device)
13         seed = generator.seed()
14
15         image = pipeline(prompt=request.prompt,
16                          negative_prompt=request.negative_prompt,
17                          generator = generator.manual_seed(seed),
18                          num_inference_steps=request.num_inference_steps,
19                          width = request.width,
20                          height = request.height,
21                          guidance_scale = request.guidance_scale
22
23                          ).images[0]
24
25         img_byte_arr = BytesIO()
26         image.save(img_byte_arr, format='PNG')
27         img_byte_arr.seek(0)
28         return StreamingResponse(img_byte_arr, media_type="image/png")
29     except Exception as e:
30         raise HTTPException(status_code=500, detail=str(e))
31
```

CLIENT API - I



FAST API'ye gönderilecek request'e user girdisinin olmadığı senaryolar

- **Styled Avatar Generation**

```
1 type ImageRequest struct
2     Prompt      string `json:"prompt"`
3     NegativePrompt string `json:"negative_prompt"`
4     NumInferenceSteps int `json:"num_inference_steps"`
5     Width        int `json:"width"`
6     Height       int `json:"height"`
7     GuidanceScale float32 `json:"guidance_scale"`
8 }
```

```
1 func main() {
2
3     requestData := ImageRequest{
4         Prompt:
5         "A potrait of a man , comic styled, detailed hair, smiling face",
6         NegativePrompt: "nsfw, bad eyes, bad teeth",
7         NumInferenceSteps: 40,
8         Width: 512,
9         Height: 512,
10        GuidanceScale: 7.5,
11    }
12
13    jsonData, err := json.Marshal(requestData)
14
15    resp, err := http.Post("http://localhost:8000/txt2img/",
16        "application/json", bytes.NewBuffer(jsonData))
17
18    defer resp.Body.Close()
19
20    outFile, err := os.Create("output.png")
21
22    _, err = io.Copy(outFile, resp.Body)
23
24    log.Println("Image saved as output.png")
25 }
```

CLIENT API - II



FAST API'ye gönderilecek
request'e user girdisi olan
senaryolar

- Stock Image Generation

```
1 func main() {  
2     http.HandleFunc("/", proxyHandler)  
3     log.Println("Starting server on :8080")  
4     log.Fatal(http.ListenAndServe(":8080", nil))  
5 }
```

```
1 func proxyHandler(w http.ResponseWriter, r *http.Request) {  
2  
3     client := &http.Client{}  
4     req, err := http.NewRequest("POST",  
5         "http://localhost:8080/txt2img/", r.Body)  
6  
7     req.Header.Set("Content-Type", "application/json")  
8  
9     resp, err := client.Do(req)  
10  
11     io.Copy(w, resp.Body)  
12 }
```

CLIENT API - III

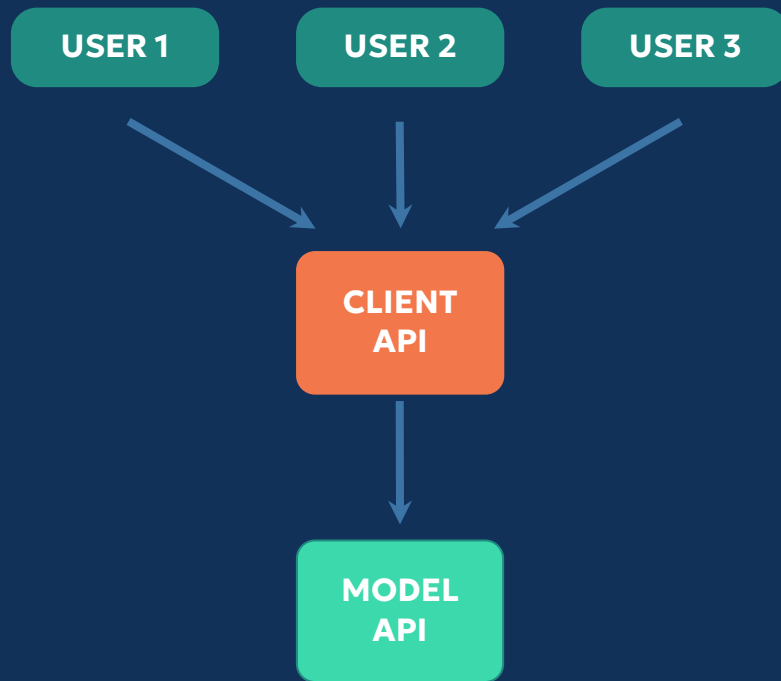


User tarafından eş zamanlı
birden fazla request geldiği
senaryolar

*Goroutines

```
1 var requestQueue = make(chan []byte, 100)
```

```
1 func main() {  
2     go processRequests()  
3     http.HandleFunc("/", proxyHandler)  
4     log.Println("Starting server on :8080")  
5     log.Fatal(http.ListenAndServe(":8080", nil))  
6 }  
7
```





- Request Body'yi okuyup byte[] kaydeder
- Eğer queue dolmadıysa değerleri Queue'ya ekler, dolduysa http error gönderir.

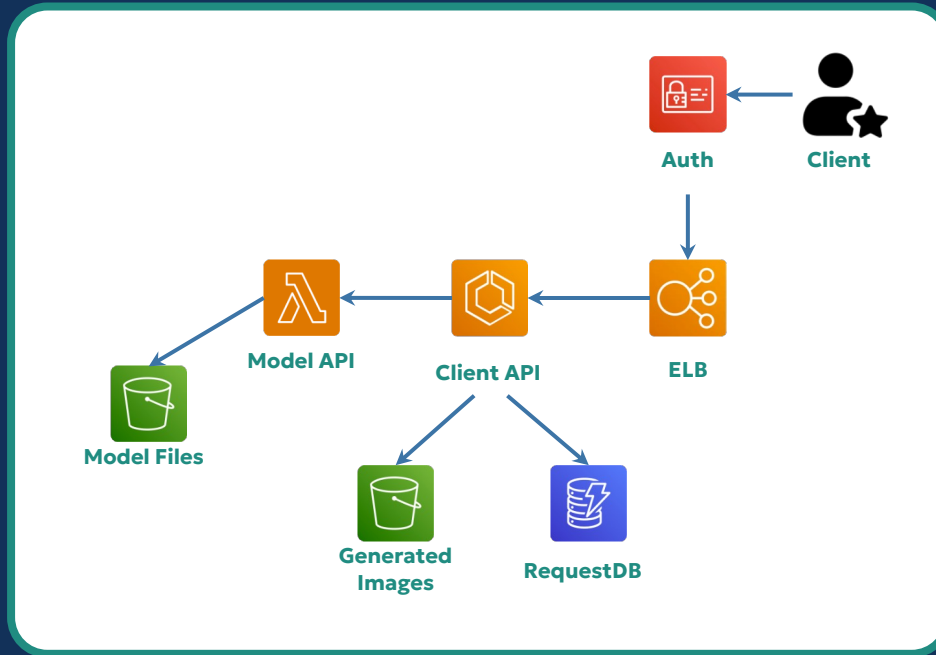
```
1 func proxyHandler(w http.ResponseWriter, r *http.Request) {  
2  
3     bodyBytes, err := ioutil.ReadAll(r.Body)  
4  
5     defer r.Body.Close()  
6  
7     select  
8     case requestQueue <- bodyBytes:  
9         log.Println("Request queued")  
10        fmt.Fprintln(w, "Request queued")  
11    default:  
12        http.Error(w, "Request queue is full", http.  
13            StatusServiceUnavailable)  
14    }
```



- Sonsuz döngü
- Queue'dan request bekler
- Request geldiği zaman bir goroutine başlatır
- Her request'i ayrı bir goroutine'de işler

```
1 func processRequests() {
2     for
3     {select
4         case bodyBytes := <-requestQueue:
5             go func(bodyBytes []byte) {
6                 client := &http.Client{}
7                 req, err := http.NewRequest("POST",
8                     "http://localhost:8000/txt2img/", bytes.NewBuffer(bodyBytes))
9                 if err != nil
10                    log.Println("Error creating request:", err)
11                    return
12                }
13                req.Header.Set("Content-Type", "application/json")
14                resp, err := client.Do(req)
15                if err != nil
16                    log.Println("Error sending request:", err)
17                    return
18                }
19                defer resp.Body.Close()
20            }(bodyBytes)
21        }
22    }
23 }
24 }
```

ÖRNEK SERVERLESS AWS MİMARİSİ



Gokonf 2024
17 Şubat
Yağmur YILDIZ

TEŞEKKÜRLER!



/in/ygryildiz



@yagmurudurdurabilirmisin



@yagmurxyildiz



@yagmurx



yagmur.cc