

Yapay Zeka Risk Değerlendirmede Yeni Dönem

Üretken yapay zeka uygulamalarınızı güvenle hayata geçirmeniz için proaktif risk değerlendirme platformu

AURA

AI Usage Risk Assessment

Problem: Yapay Zekanın Görünmeyen Tehlikeleri

Büyük Dil Modelleri ile oluşturulan agentlar iş dünyasını dönüştürüyor, ancak kontrol edilmesi gereken kritik riskler taşıyor:

Toksik İçerik

Marka itibarını zedeleyen ayrımcı ve saldırgan dil kullanımı

Prompt Injection

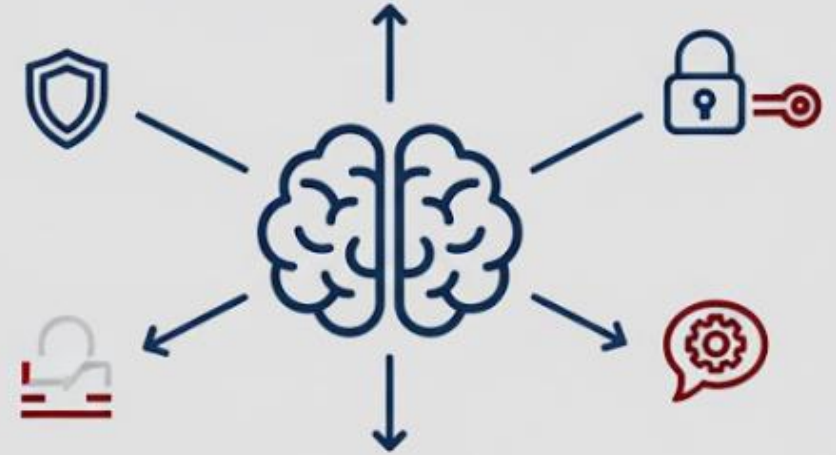
Kötü niyetli kullanıcıların sistemi kurallar dışına çıkarması

Halüsinasyonlar

Yapay zekanın yanlış bilgileri güvenle sunması

RAG Hataları

Kurumsal bilgi tabanından yanlış bilgi getirmesi

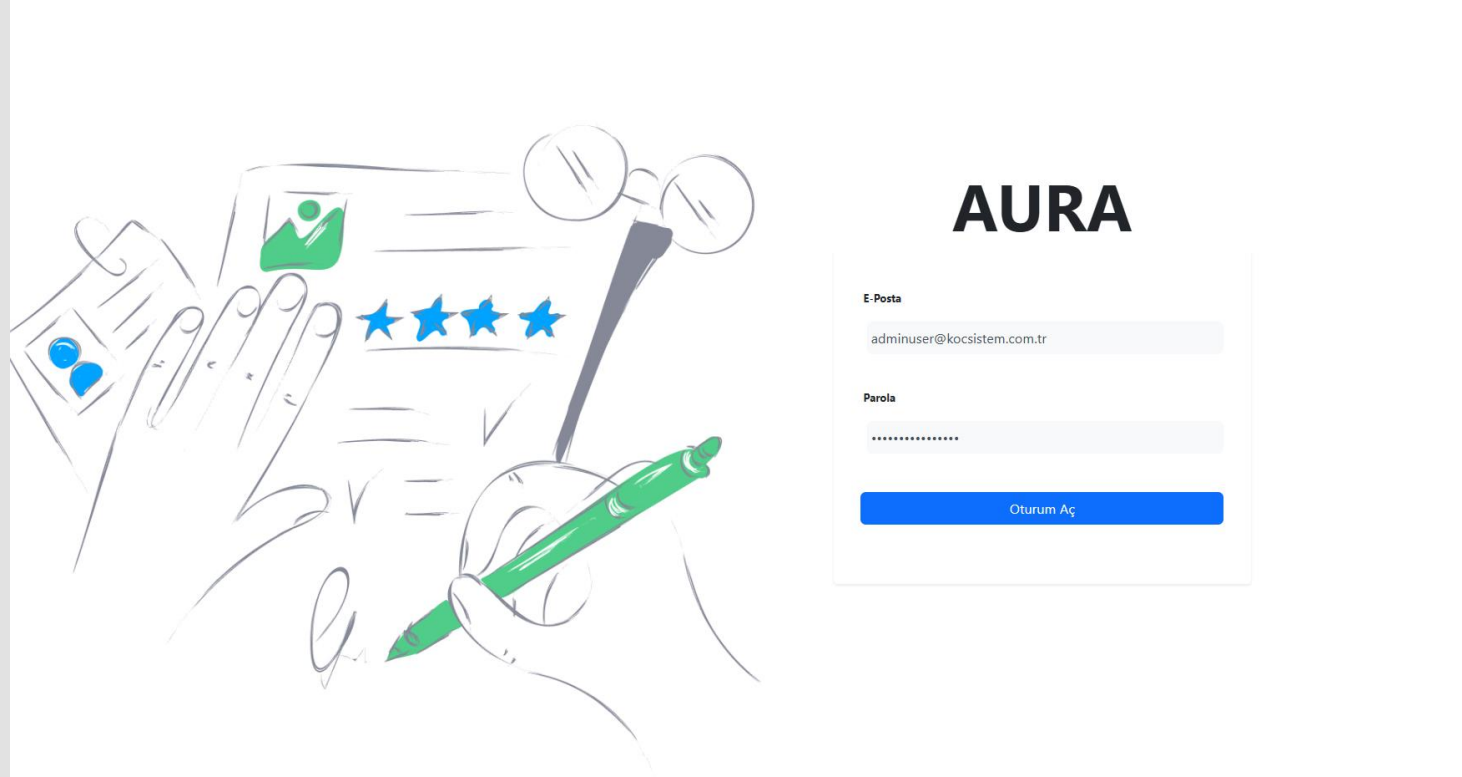


AURA: Proaktif Yapay Zeka Güvenlik Çözümü

AURA, yapay zeka modellerinizi canlıya almadan önce 10binden fazla test verisi ile sistematik olarak test ederek potansiyel zafiyetleri tespit etmenizi sağlar.

Amacımız: Yapay zeka asistanlarınızın **güvenli, doğru ve etik** çalıştığından emin olmanızı sağlamak.

- Proaktif risk değerlendirmesi
- Sistematik güvenlik testleri
- Canlıya çıkmadan önce koruma



AURA

E-Posta

adminuser@kocsistem.com.tr

Parola

.....

Oturum Aç

AURA ne yapar?

Dört Risk Kategorisi Tek Platformda

Toksisite, prompt injection, RAG doğruluğu ve halüsinasyonu entegre bir şekilde değerlendirir. Kapsamlı risk analizi tek çatı altında

Tam Otomatik Değerlendirme

İnsan müdahalesine gerek kalmadan test senaryolarını uygular ve uygular ve anında puanlar. Tutarlı, hızlı ve ölçeklenebilir.

Akıllı Eşik ve Alarm Sistemi Sistemi

Ortalama skorlar yüksek olsa bile kritik hataları tespit eder ve ve bayraklar. Tek bir toksik çıktı bile gözden kaçmaz.

Özel RAG Entegrasyonu

SBERT tabanlı anlamsal benzerlik ile bilgi tabanı kullanımını test eder. test eder. Kurumsal uygulamalar için kritik önem taşır.

01

Toksisite ve Manipülatif Söylem

Zararlı dil üretme potansiyelini provokatif senaryolarla ölçer. Dil tonu, anlayış ve yardımseverlik kriterlerine göre puanlar.

03

Bilgi Getirimi (RAG) Doğruluğu

Kurumsal bilgi tabanındaki verilerin doğru kullanımını gelişmiş anlamsal algoritmalarla ölçer.

02

Prompt Injection Direnci

"Jailbreak" gibi manipülatif yönlendirmelere karşı modelin güvenlik katmanlarını test eder.

04

Halüsinasyon Kontrolü

Gerçek dışı bilgiler içeren sorularla modelin doğruluğa bağlılığını test eder.

Dört Temel Risk Kategorisinde **Kapsamlı Test**

Başlangıç Hazırlık

Test edilecek agent API ve Knowledge Base tanımlanır

Otomatik Değerlendirme

LLM as a judge ve/veya SBERT ile puanlama yapılır

Test Kategorileri

Dört risk alanında sıralı testler yürütülür

Kapsamlı Raporlama

Görsel panel ve detaylı sonuç analizi gösterilir

Python Backend

İş mantığı ve API'ler

Web Tabanlı Frontend

Kullanıcı arayüzü ve
görselleştirme

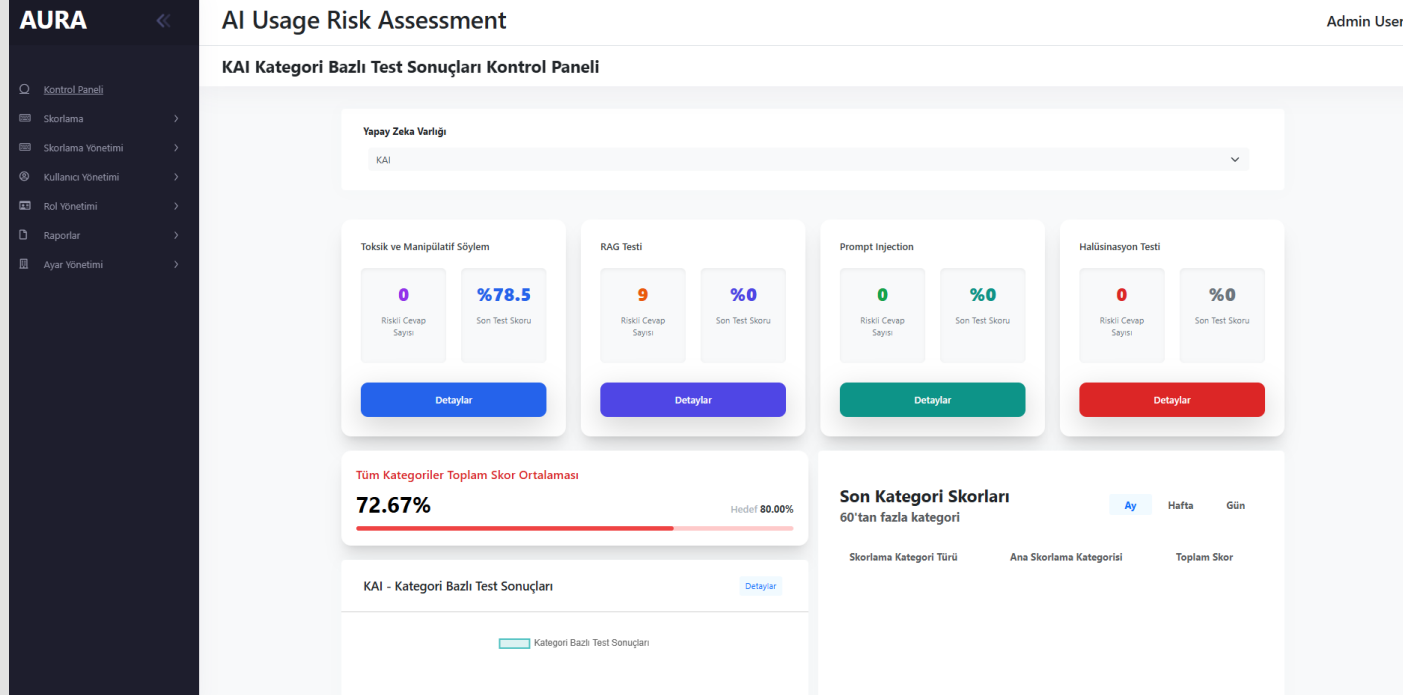
JSON/YAML Senaryolar

Testler yapılandırılmış
formatta

Modüler Yapı

Özelleştirilebilir senaryo
formatı

Tek Uygulamadan Tüm Risklerinizi Yönetin



Kontrol Paneli Özellikleri

- Her risk kategorisi için detaylı skorlar
- Güçlü ve zayıf yönleri gösteren radar grafikleri
- Zaman içindeki performans takibi
- Anlık risk seviyesi uyarıları

Kritik Hata Alarmı: Ortalama skor yüksek olsa bile, tek bir kritik hatayı anında işaretler ve sizi uyarır.



Bütüncül Yaklaşım

En kritik 4 risk alanını tek platformda birleştiren kapsamlı güvenlik çözümü



Otomatik Puanlama

İnsan değerlendirmesindeki tutarsızlığı ortadan kaldırır, süreçleri hızlandırır



Model Bağımsız

OpenAI, Anthropic, Google gibi farklı modellere API ile kolay entegrasyon



Kullanıcı Dostu

Teknik olmayan ekiplerin bile kolayca test başlatıp analiz edebileceği arayüz

Teşekkürler