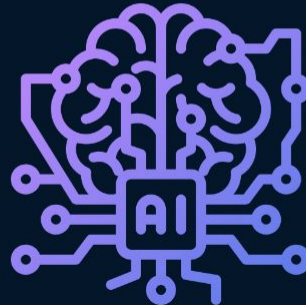


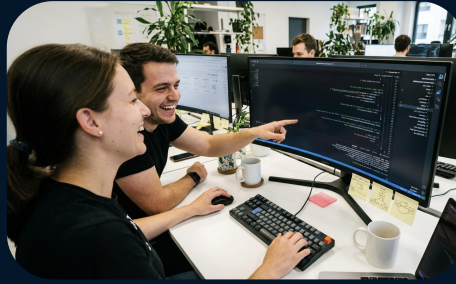
# METRO digital

## Chatbot Evaluation

*“You can’t manage, what you can’t measure”*



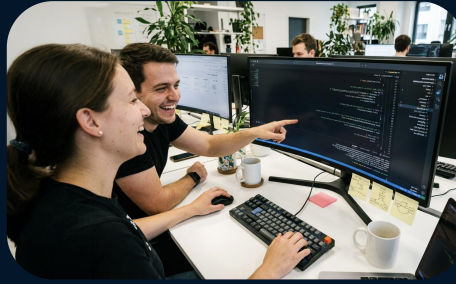
# *Problem*



AI Developer

Unser neuer FAQ  
Agent soll live gehen

# Problem



AI Developer

Danke... war  
etwas besser!

Neue Version,  
bitte testen!

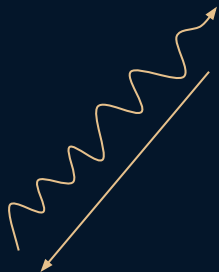


Business Unit

# Problem



AI Developer



Business Unit

- ! Manuelle Eingabe der Testfragen
- ! Kein hilfreiches Feedback
- ! Keine Benchmarks



Go life approved?



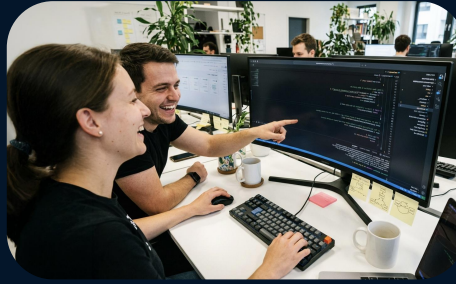
Was heißt  
besser?

Bloß  
keine  
Skandale!

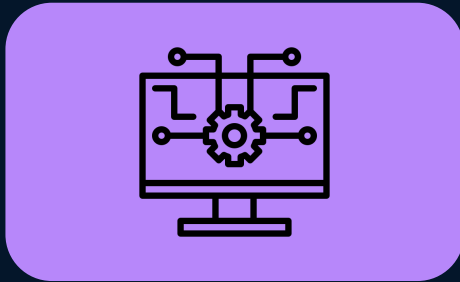


Board

# Lösung



AI Developer



AI Evaluation Tool

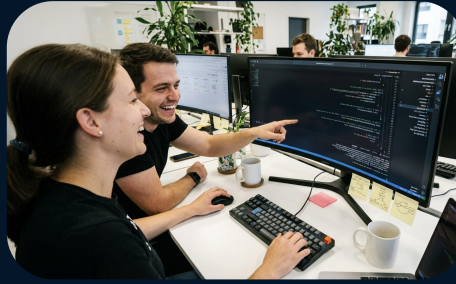
- 😍 Automatisiertes Testing
- 😍 Detaillierte Analyse
- 😍 Messbare KPIs

Go live  
approved?



Board

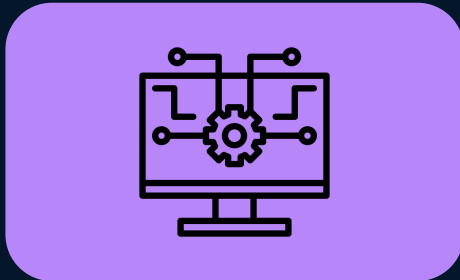
# Lösung



AI Developer



Go live  
approved?



AI Evaluation Tool

- 🥰 Automatisiertes Testing
- 🥰 Detaillierte Analyse
- 🥰 Messbare KPIs



Board

# ***LIVE Demo***

This is where the magic happens ✨

# *Pipeline*



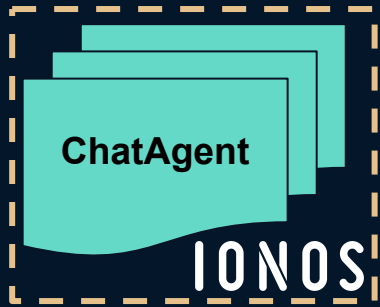
python script

csv file

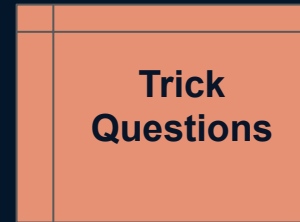
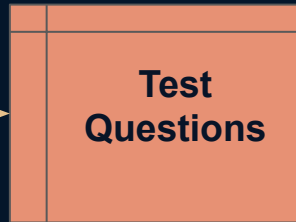
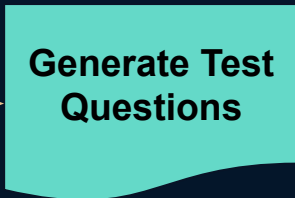
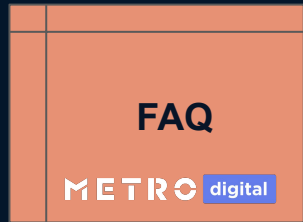
output



# Pipeline



↑ RAG



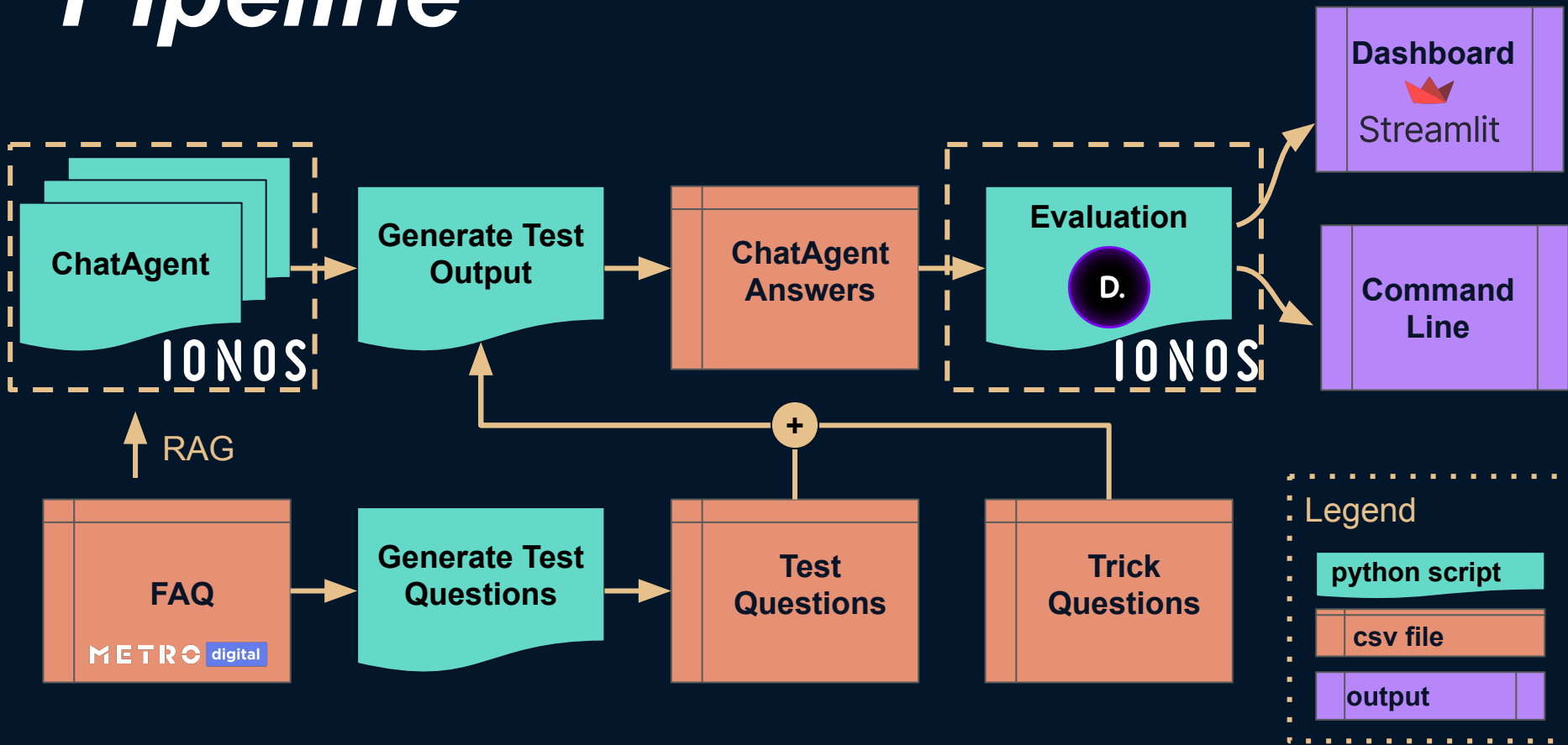
## Legend

python script

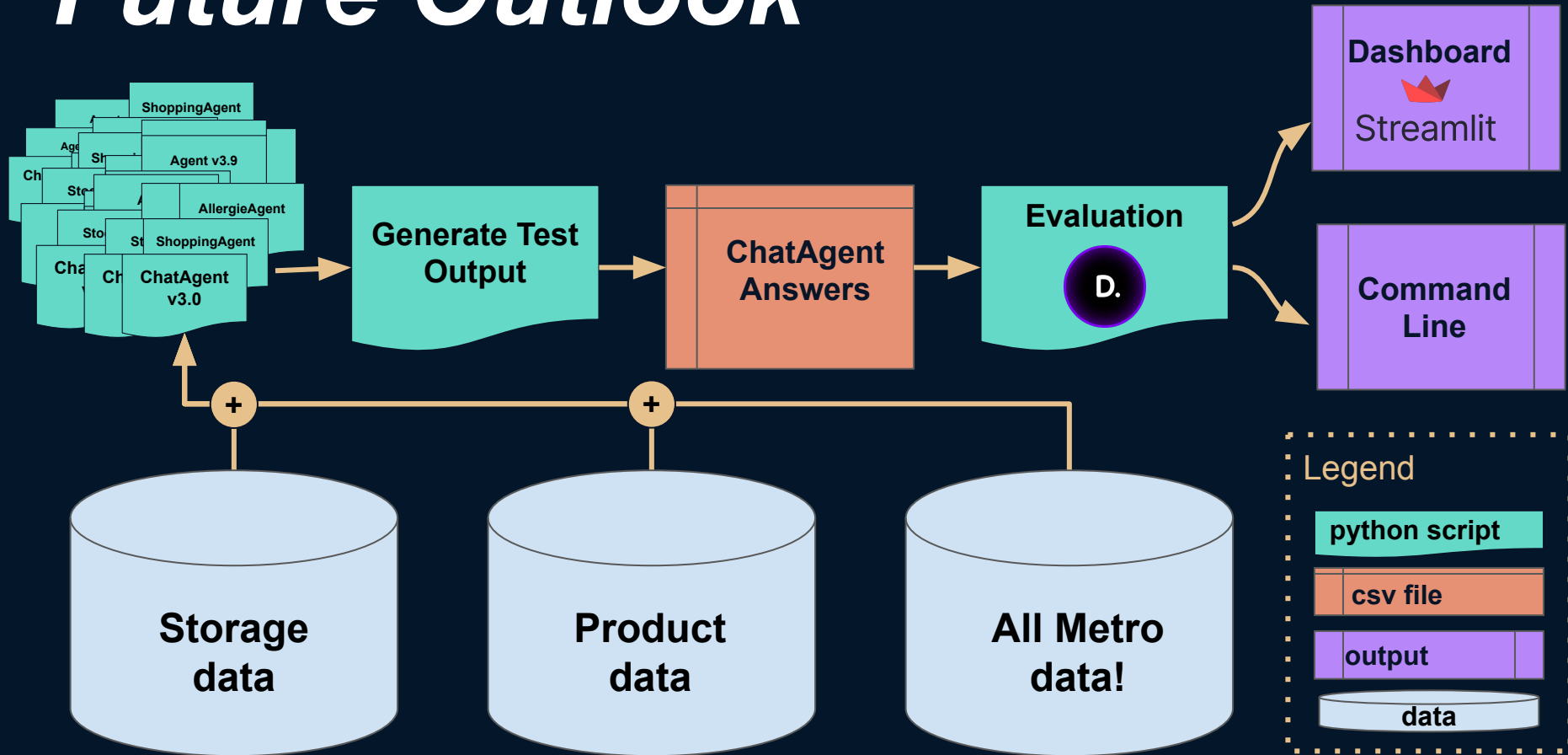
csv file

output

# Pipeline



# Future Outlook



# A success story

## METRO digital & IONOS

Manage What You Measure.



**Hauke  
Diers**



**Domitille  
Grandjean**



**Nico  
Neubauer**



**Andres  
Navarro**



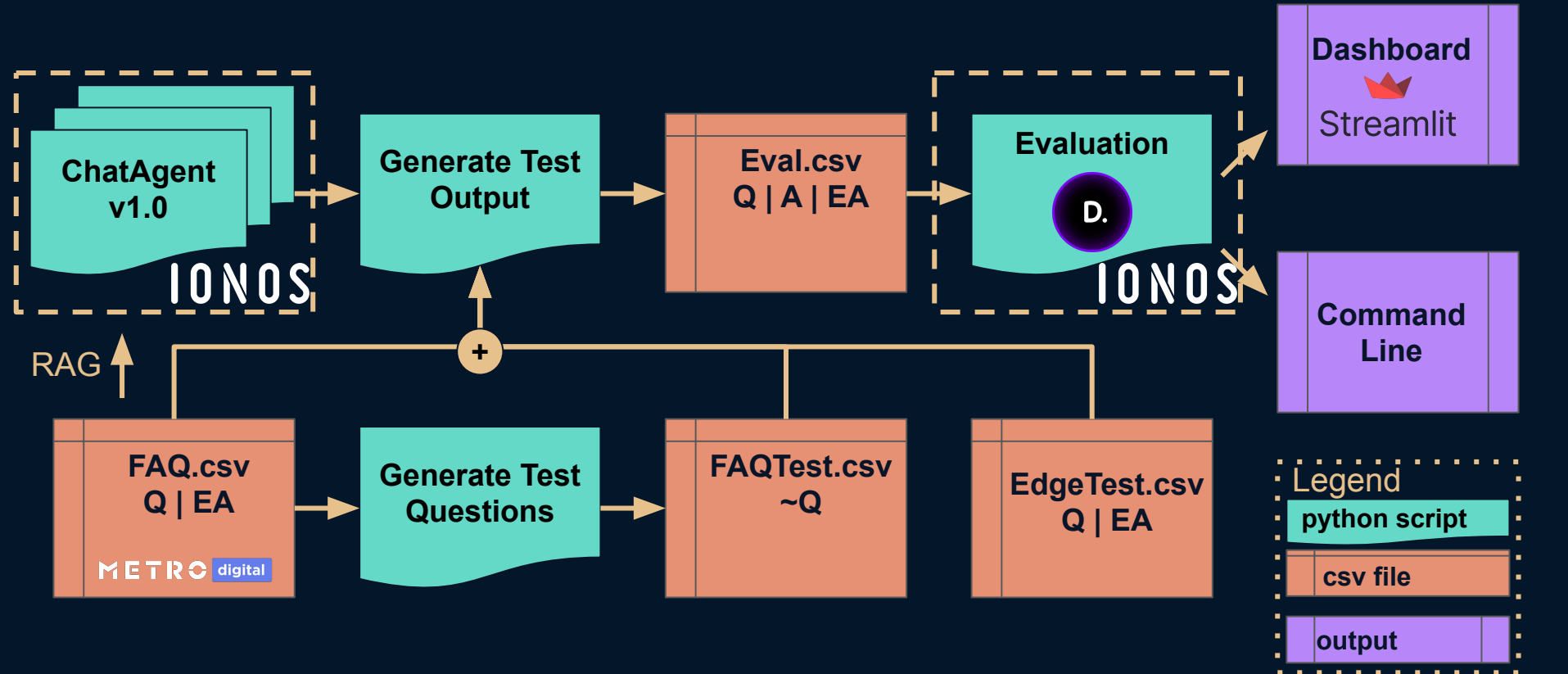
**Omid  
Roshenas**

# *Evaluation Metrics*

1. **Correctness:** Sind die angegebenen Informationen wahr?
2. **Clarity:** Ist die Antwort leicht zu verstehen?
3. **Hospitality & Tonality:** Ist der Ton professionell?
4. **Relevance:** Beantwortet er die gestellte Frage?
5. **Hallucination Safety:** Werden Informationen erfunden?



# Future Outlook



# *Skalierbarkeit*

Naheliegende Use cases:  
Dish:

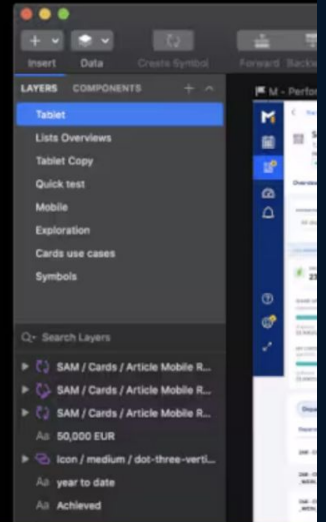
Broader use cases für weitere Integration:  
Evaluation Robotics AI - Accuracy of order  
preparation  
Maintenance check-up for Obi's new storage  
tracker

```

~\code\docs\rebase\test\index.js
DarkMode.js index.js

1 import React from 'react'
2 import { ThemeProvider } from 'styled-components'
3 import { Box } from 'grid-styled'
4
5 import { invertLuminance } from './colors'
6 import defaultTheme from './theme'
7
8 export const invertTheme = (theme = defaultTheme) => {
9   const { colors = {}, ...rest } = theme
10   const next = Object.keys(colors)
11     .reduce((a, key) => {
12       a[key] = invertLuminance(colors[key])
13       return a
14     }, {})
15
16   return {
17     ...rest,
18     colors: next
19   }
20 }
21
22 export class DarkMode extends React.Component {
23   static defaultProps = {
24     color: 'black',
25     bg: 'white'
26   }
27
28   render() {
29     return (
30
31
32   test('Provider renders with custom the
33   const json = render(<Provider
34     theme={
35       fonts: [],
36       fontSizes: [
37         12, 16, 18, 24, 36, 48, 72
38       ],
39       space: [
40         0, 6, 12, 18, 24, 30, 36
41       ]
42     })
43   />).toJSON()
44   expect(json).toMatchSnapshot()
45 })
46
47 test('theme is an object', () => {
48   expect(typeof theme).toBe('object')
49   expect(Array.isArray(theme.breakpoints)).t
50   expect(Array.isArray(theme.space)).t
51   expect(Array.isArray(theme.fontSizes)).t
52   expect(typeof theme.fontWeights).toBe('ob
53   expect(typeof theme.colors).toBe('obj
54   expect(typeof theme.radii).toBe('obj
55   expect(typeof theme.fonts).toBe('obj
56 })
57
58 examples.forEach(( { name, element } ) => {
59   test(`${name} renders`, () => {
60     const { container } = renderWithTheme

```





# *Welche Metriken Evaluation*

Correctness

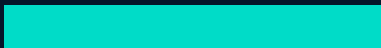
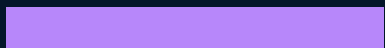
Clarity

Hospitality / Tonality

Relevance

Hallucination Safety

Each scale 1/5



# Pipeline

