



A origem da complexidade - Dívida técnica na área de dados

André Gomes dos Santos

Diamante



Platina



DISCOVER

Ouro



VERTICA
by opentext

Prata

TRACES



Apoio

FIA.P



Who am I ?

<https://www.linkedin.com/in/andregoems/>



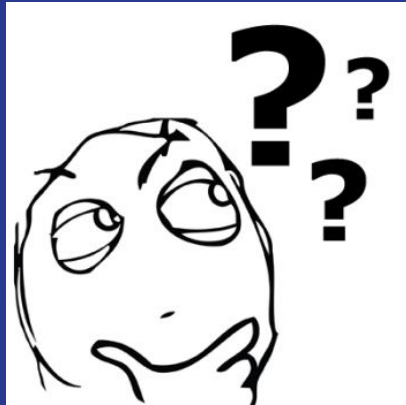
!

●

Disclaimer



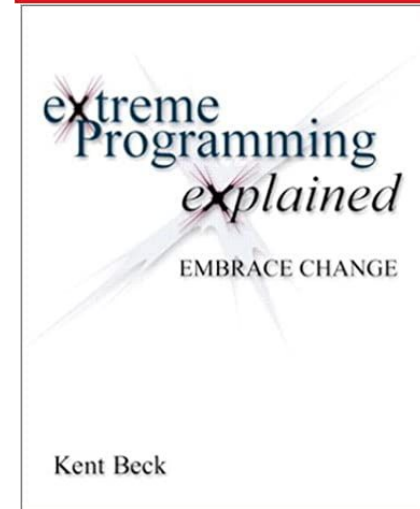
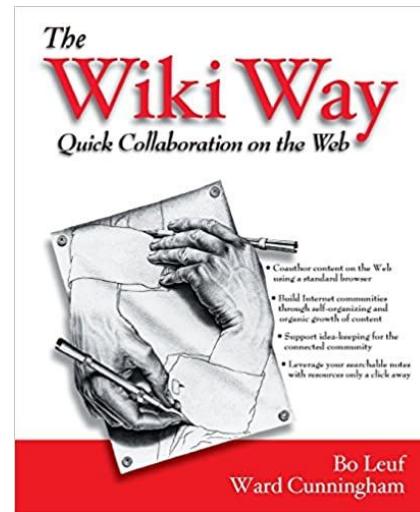
Dívida (débito) técnica ?



Ward Cunningham



É uma **metáfora** criada para explicar a **natureza evolutiva** de um software, e também o **custo de aprendizado e refatoração** no decorrer da vida da aplicação.

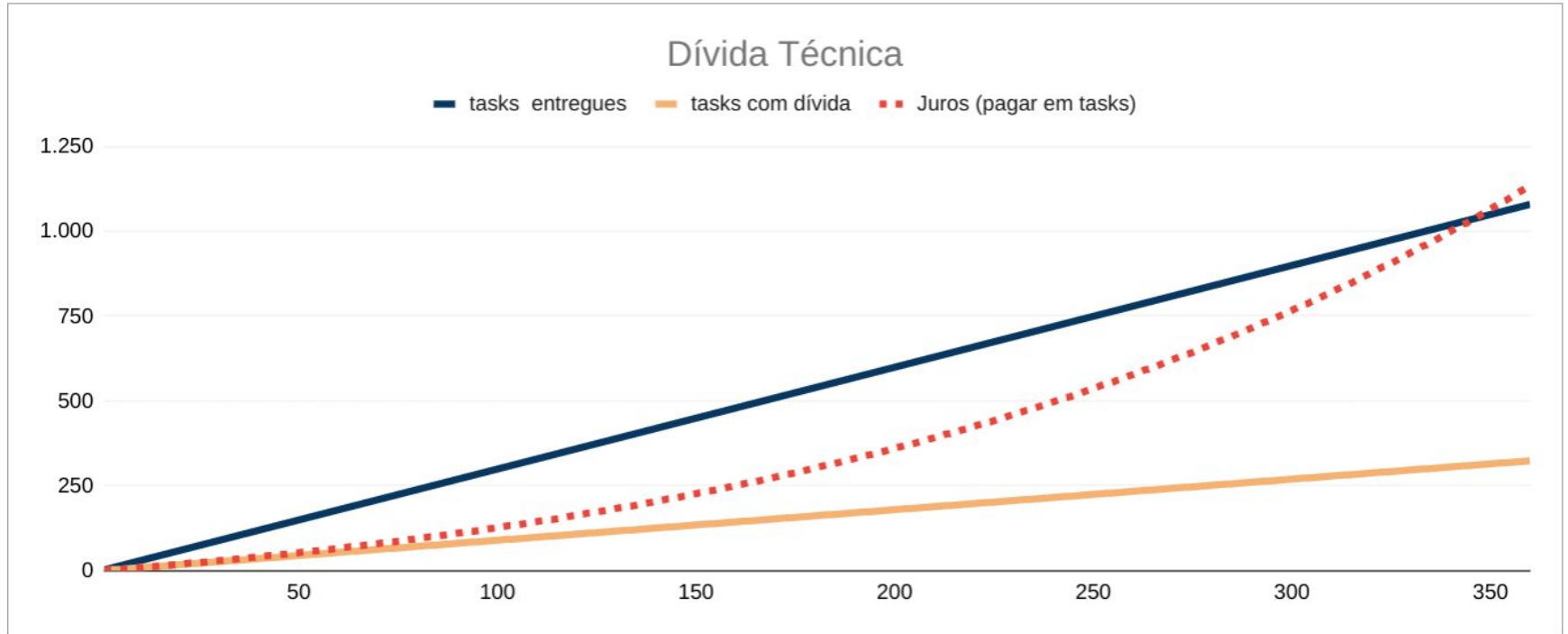




Evento	Valores
Total pago no financiamento	R\$ 216.797,75
Valor financiado	R\$ 100.000,00
Juros devidos	R\$ 116.797,75
Pagamento mensal estimado	R\$ 3.613,30

360 dias

1.080 tasks (esforço/custo/tempo/time)



Faturamento de 50.000 R\$
em dias normais

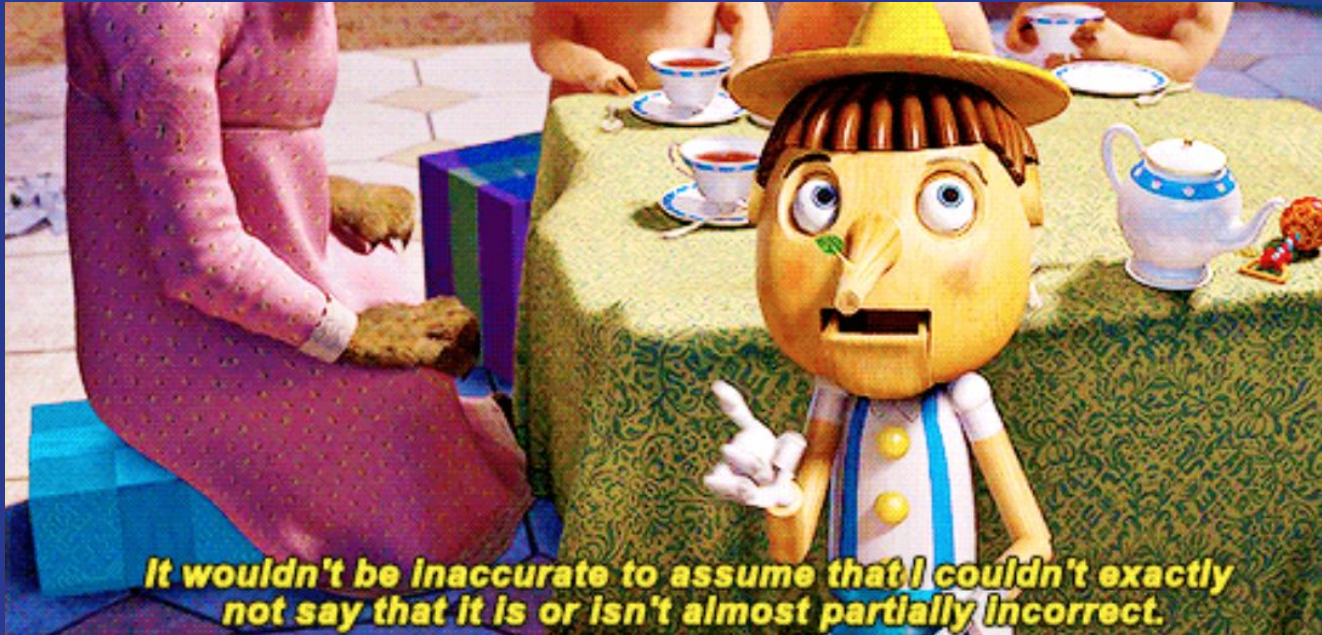


Faturamento de 300.000 R\$
em feriados.

Evento	Valores
Total pago no financiamento	R\$ 216.797,75
Valor financiado	R\$ 100.000,00
Juros devidos	R\$ 116.797,75
Pagamento mensal estimado	R\$ 3.613,30

- **task/job/feature/tabela/relatório entregue/feito/desenvolvido sem valor percebido é custo.**
- **Dívida deve ser controlada.**
 - Dívida deve ser adquirida conscientemente.
 - Dívida deve ser paga.

Falácias



As 8 falácias da computação distribuída.

1. A rede é de confiança.
2. A latência é zero.
3. Largura de banda é infinita.
4. A rede é segura.
5. A topologia não muda.
6. Há somente um administrador.
7. O custo de transporte é zero.
8. A rede é homogênea.

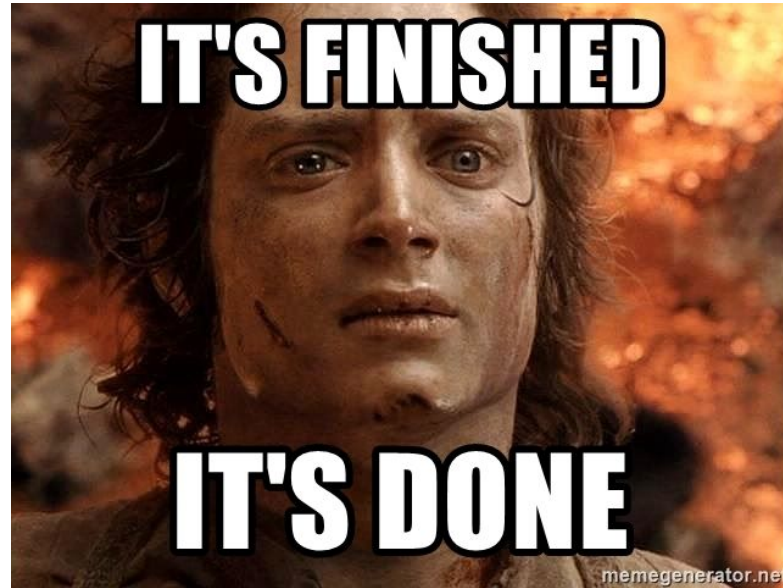
<https://cybernetus.com/post/as-oito-falacias-da-computacao-distribuida/>

Falácias da computação empresarial. (Grandes projetos/empresas)

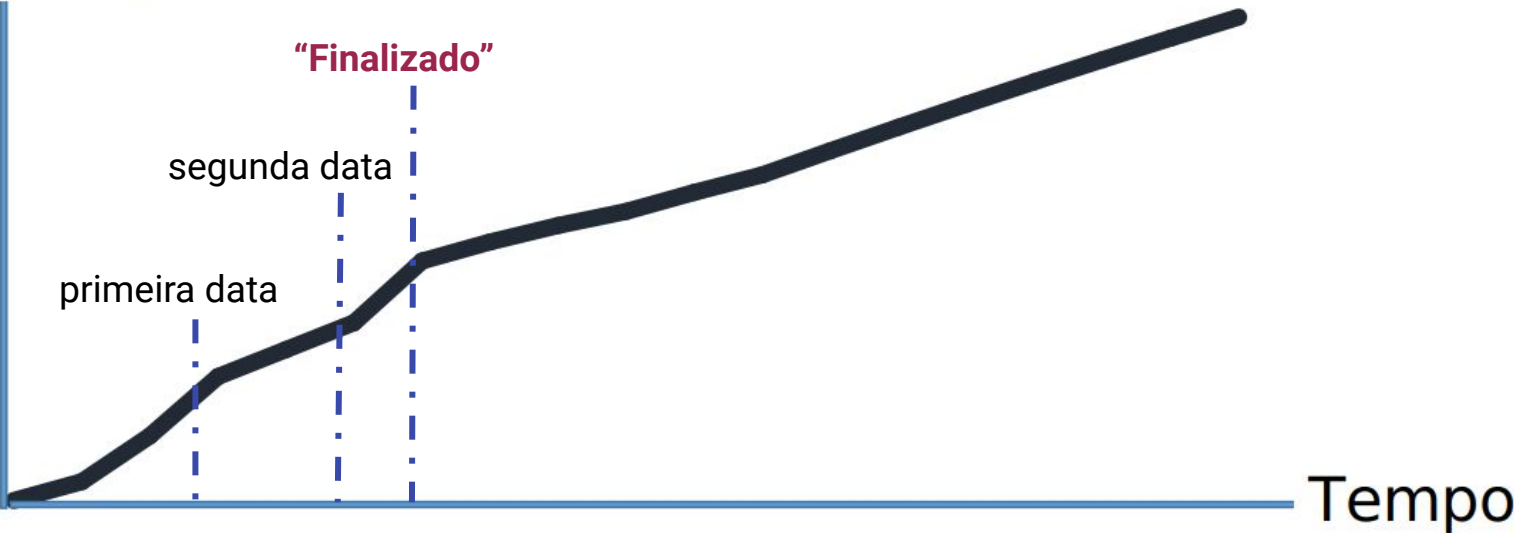
1. **A nova tecnologia é sempre melhor do que a tecnologia antiga.**
2. Sistemas não são “sistemas distribuídos”.
3. **A lógica de negócios pode e deve ser centralizada.**
4. **Dados, objetos ou qualquer outro tipo de modelo podem ser centralizados.**
5. O sistema é monolítico.
6. **O sistema está finalizado.**
7. Os fornecedores podem fazer com que os problemas desapareçam.
8. **A arquitetura é a mesma em todos os lugares.**
9. Os desenvolvedores só precisam se preocupar com problemas de desenvolvimento.

<http://blogs.tedneward.com/post/enterprise-computing-fallacies/>

Falácia: O sistema está finalizado.

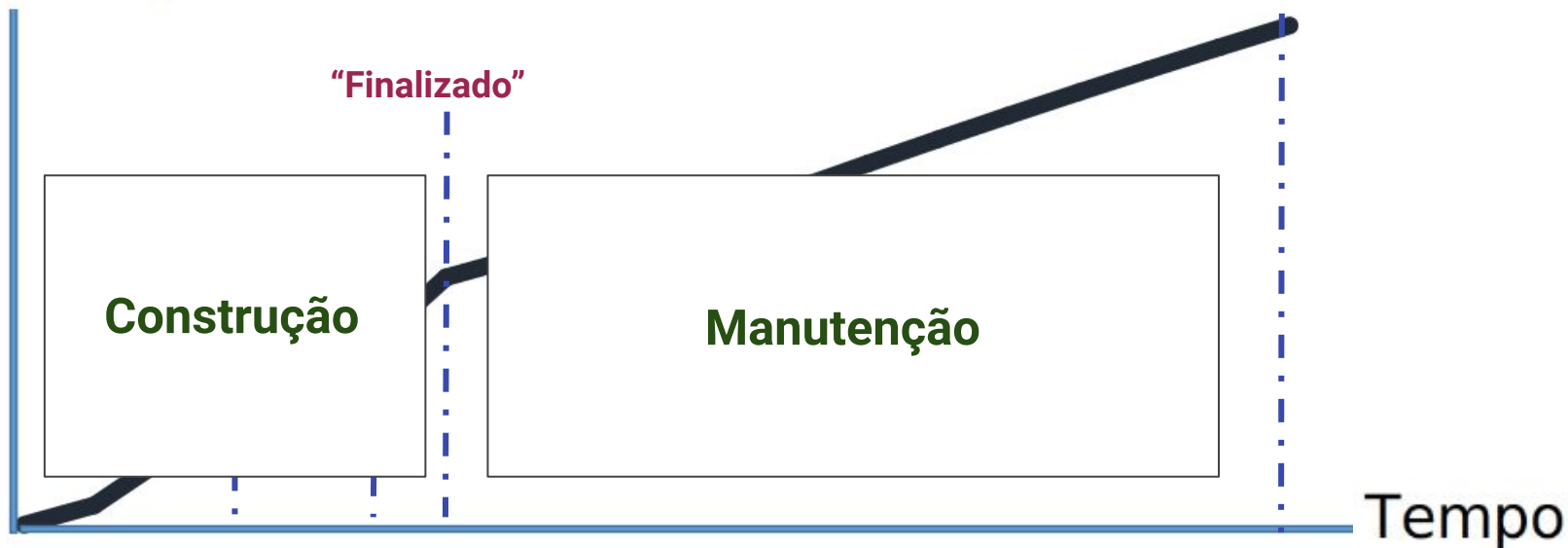


Custo R\$



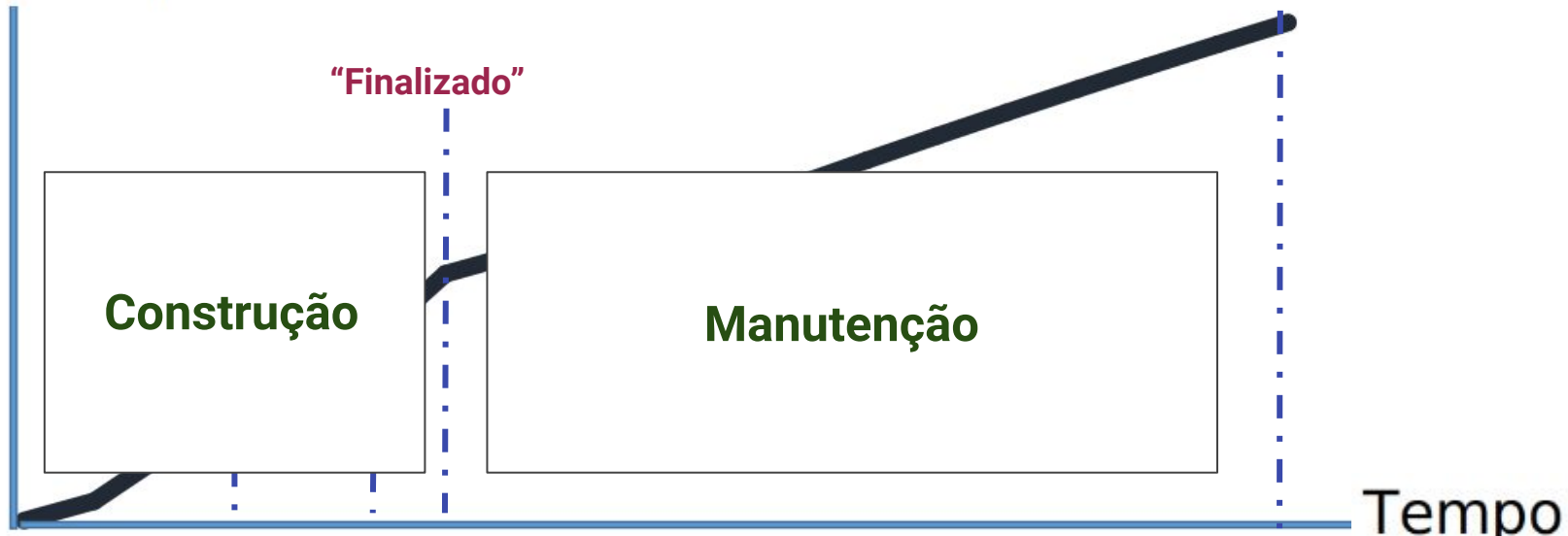
Tempo

Custo R\$



Tempo

Custo R\$



"Finalizado"

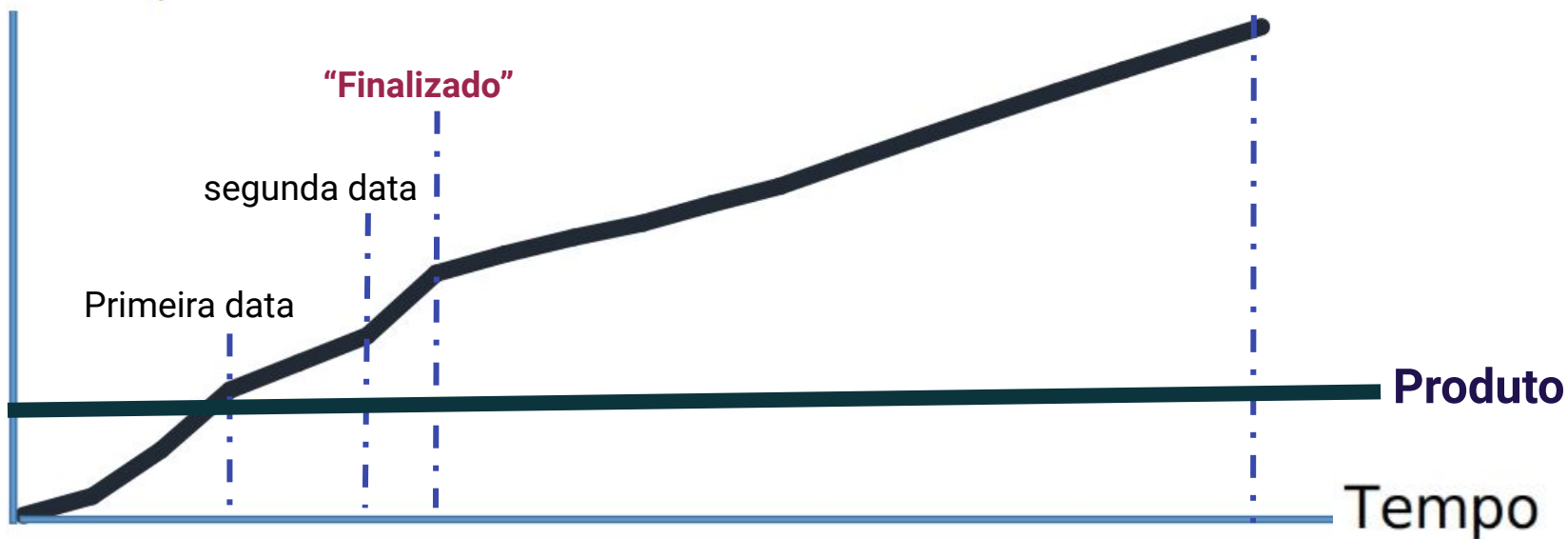
Construção

Manutenção

Temos que refazer...

Tempo

Custo R\$

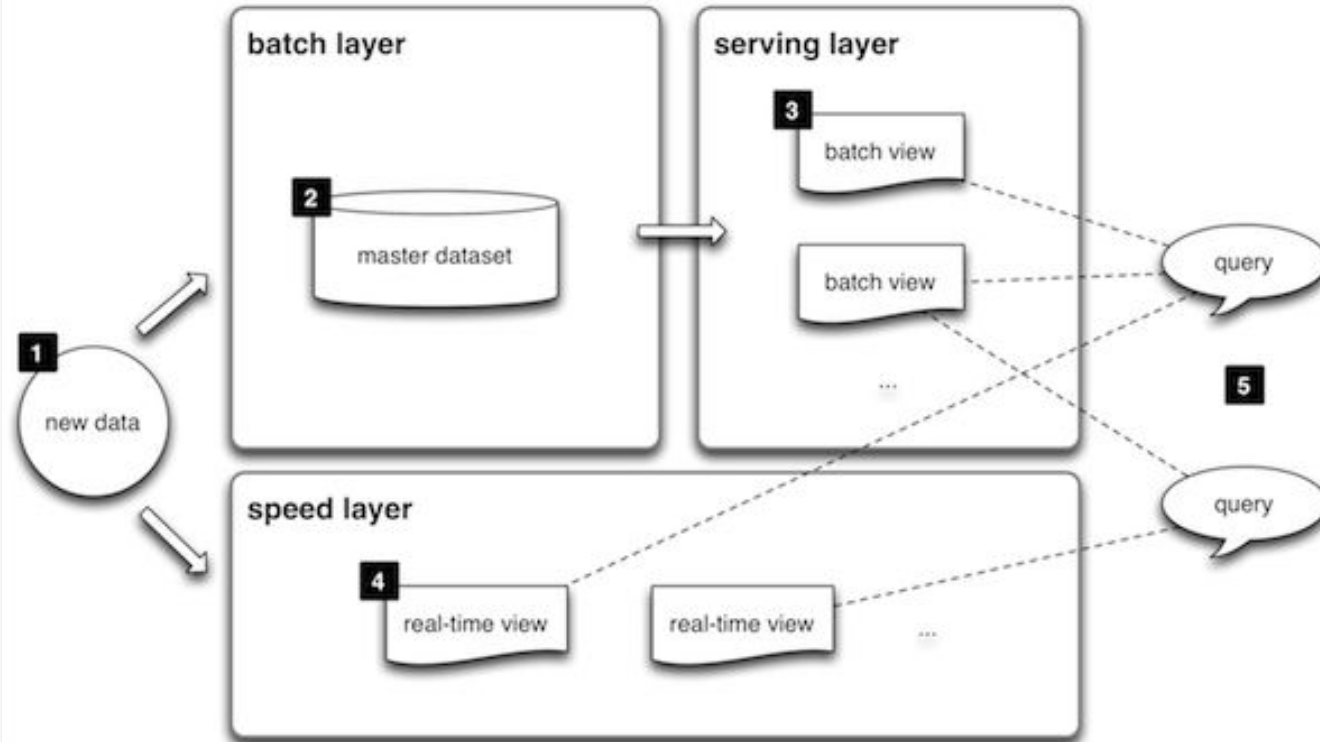


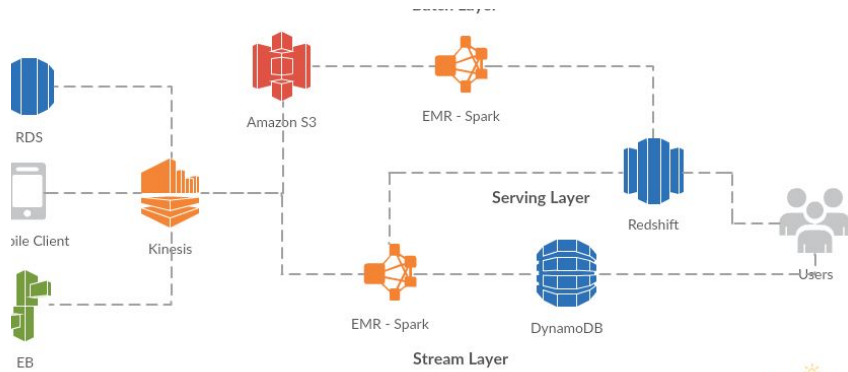
- **Desenvolvimento (dados) não é um projeto de construção civil, a metáfora induz a “erros”.**
- Não existe um “momento de manutenção”.
- Refazer tudo não é a solução (eventualmente é necessário).
- Iniciativas em **TI só terminam quando o “código” é removido/descartado.**
- Iniciativas em **Dados só terminam quando o “dado” é removido/descartado.**
- **Projetos/Iniciativa = Produtos:**
 - **Precisa de feedback do usuário.**
 - **Precisa gerar/demonstrar valor/utilidade.**
 - **é uma iniciativa de longo período.**
 - **<https://martinfowler.com/articles/products-over-projects.html>**

Falácia: A arquitetura é a mesma em todos os lugares.



Arquitetura Lambda





[online diagramming & design] [createfy.com](https://www.createfy.com)

Na AWS

Data Lake

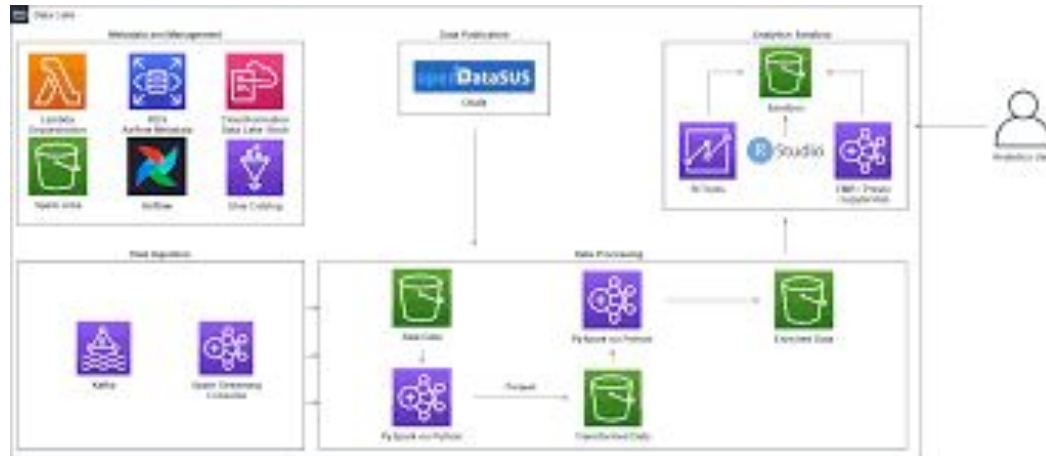
- Armazenamento: **S3**
- Metadata: **Glue**

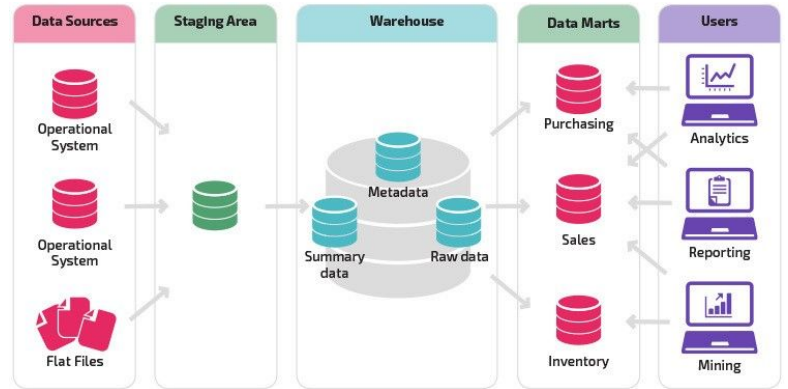
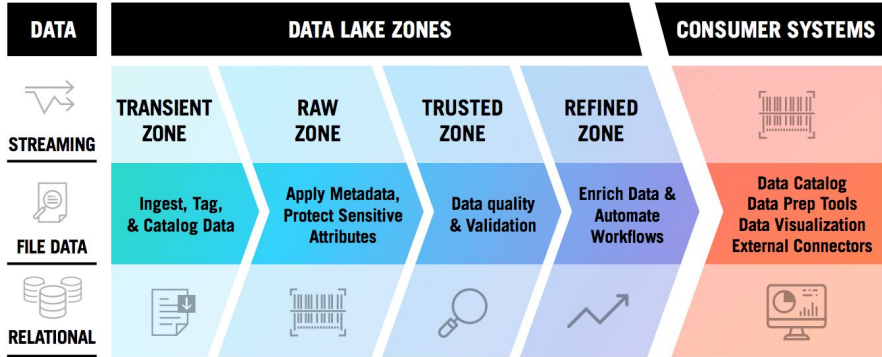
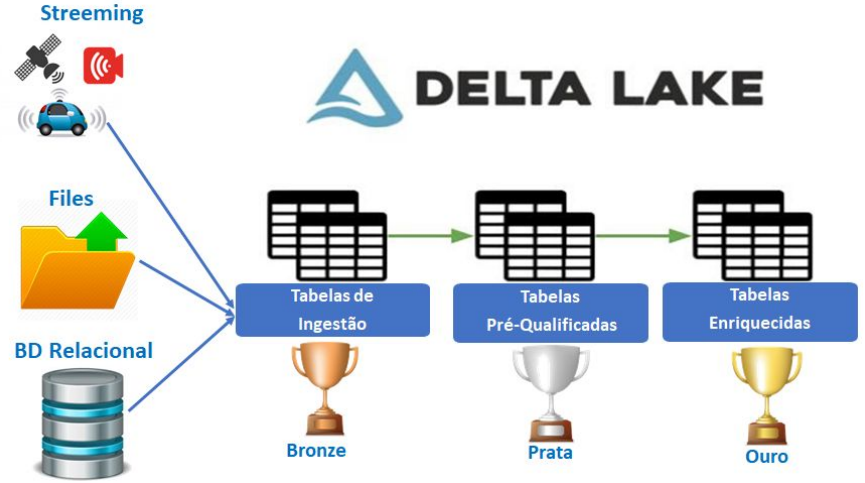
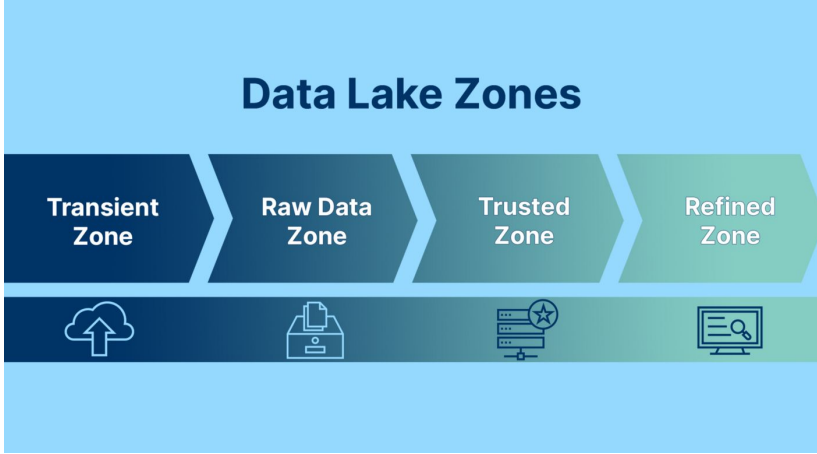
Lambda

- Processamento: **Kinesis**, Kafka, Spark, etc

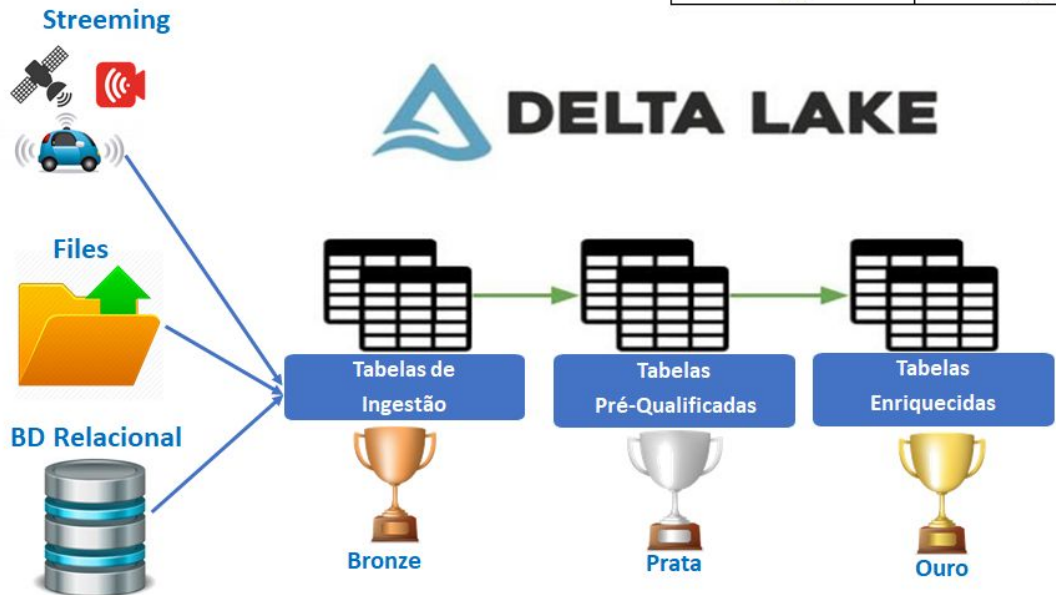
Serving

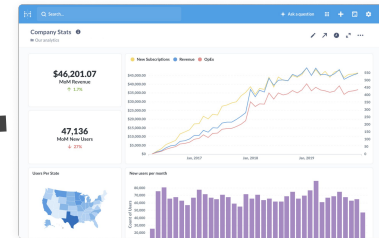
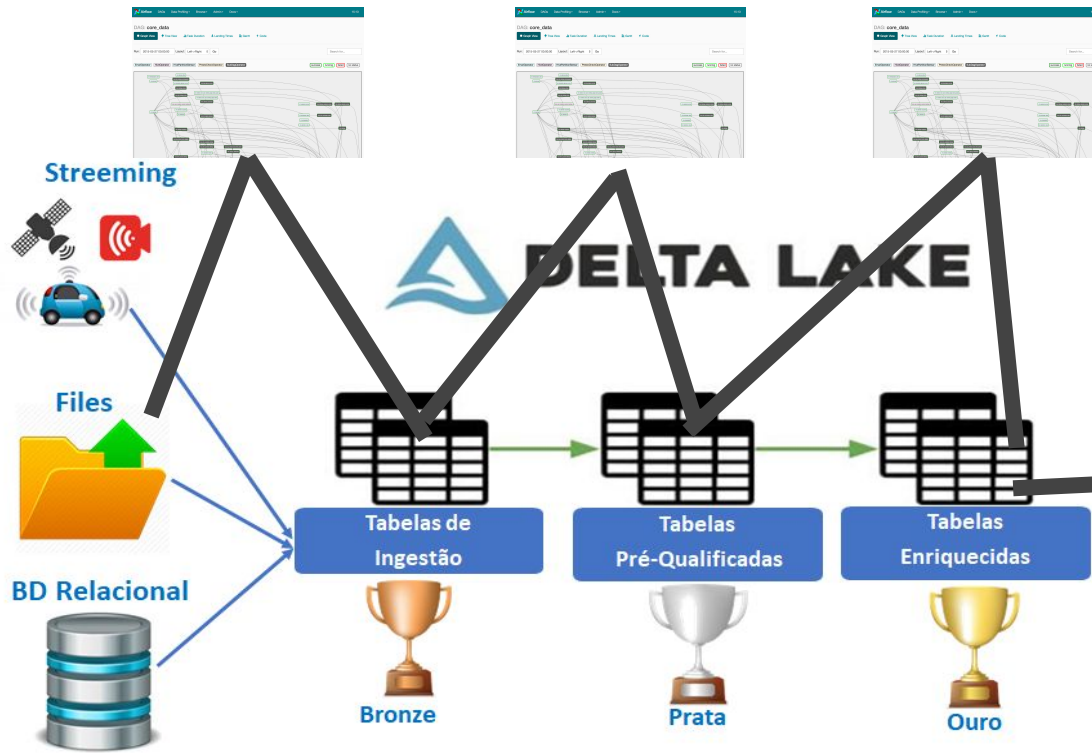
- Data Warehouse: Redshift
- Machine Learning: Sagemaker
- Analytics: Quicksight, Elastic Spectrum, Presto, etc
- Ad-Hoc: **Athena**, Redshift

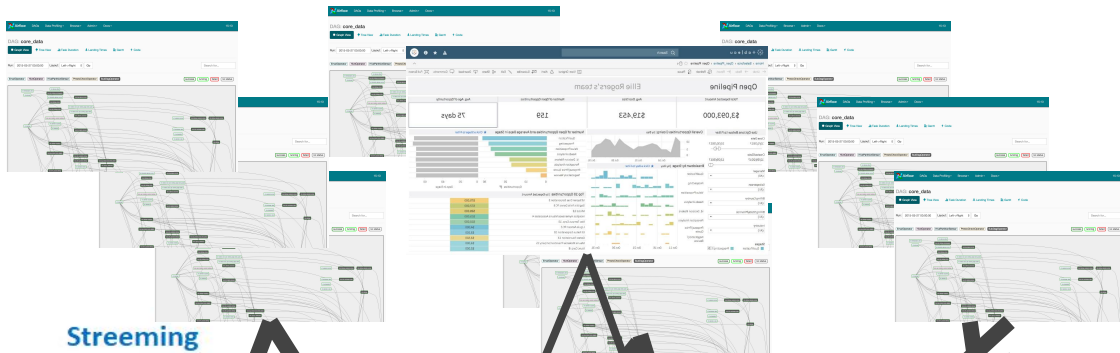




Tier1	Tier2	Tier3
raw	trusted	consumption
raw	trusted	refined
Bronze	Silver	gold
Landing	Curated	production
Raw	Conformed	Modeled
land	Prepare	publish
....







Streaming



Files



BD Relacional



DELTA LAKE



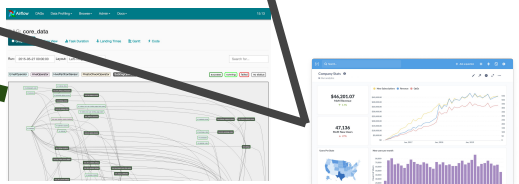
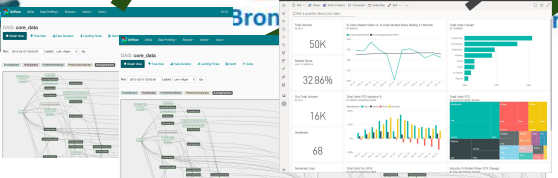
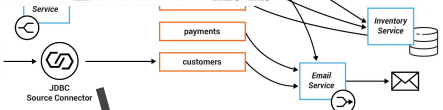
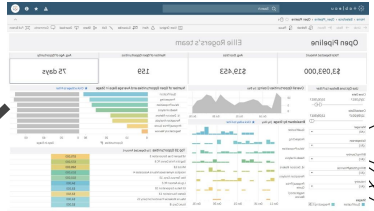
Bronze



Prata



Ouro



Streaming



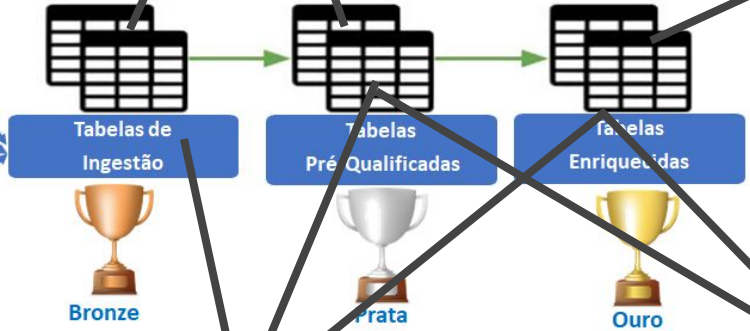
Files



BD Relacional



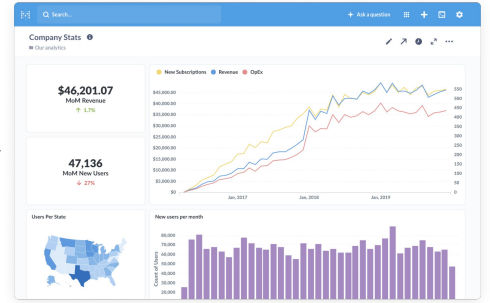
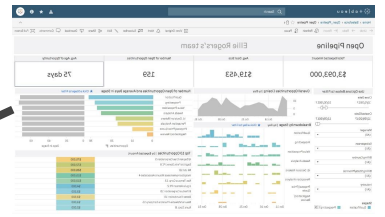
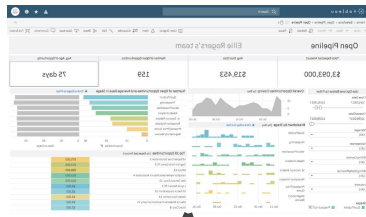
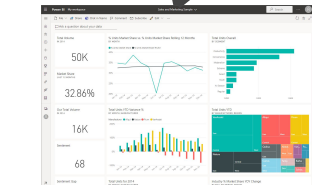
DELTA LAKE



Bronze

Prata

Ouro



Streaming



Files



BD Relacional



DELTA LAKE



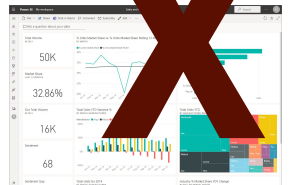
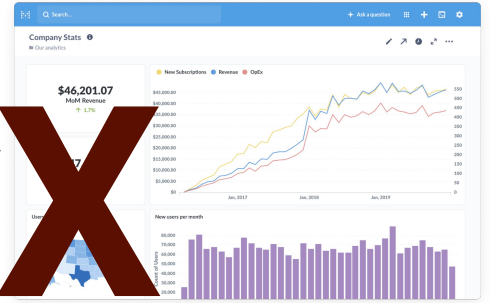
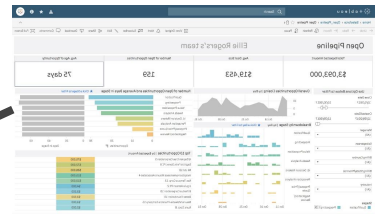
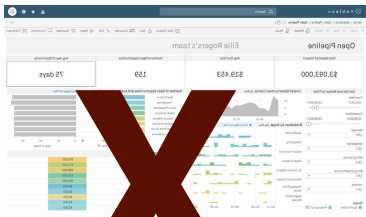
Bronze



Prata



Ouro





DELTA LAKE

Streaming



Files



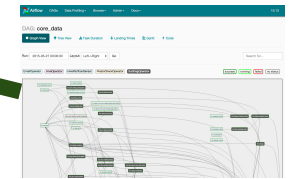
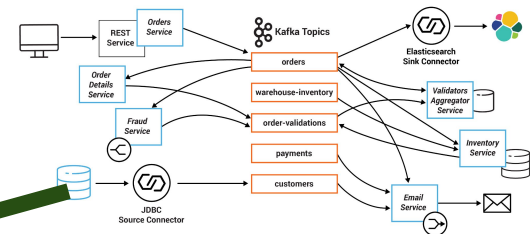
BD Relacional



Bronze

Prata

Ouro





Streaming



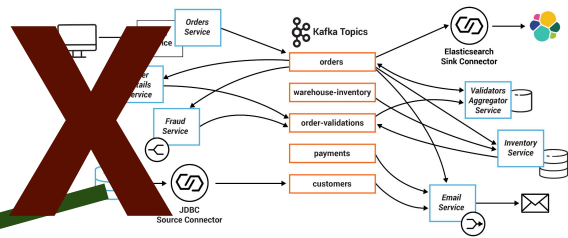
Files



BD Relacional



DELTA LAKE



Disponibilidade(escala):

Backend = nanosegundo/milisegundo/segundo

Datalake = minuto/hora/dia

Acomplamento:

depende do Datalake, mas logicamente depende de outro backend



Streaming



Files



BD Relacional



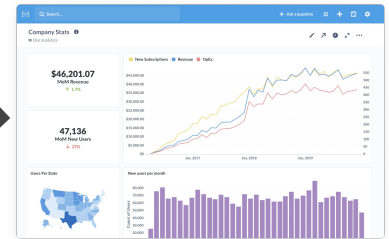
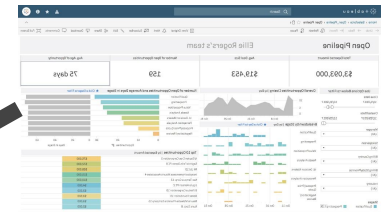
Bronze



Prata



Ouro



- **A restrição(objetiva) é o que permite a longevidade da arquitetura.**
- **A restrição(objetiva) sustenta a flexibilidade.**
- **é sempre contextual.**
- **A arquitetura “correta” é a mais “simples” de mudar (Refatorar).**
- **complexidade não deixa de existir, apenas muda de lugar.**

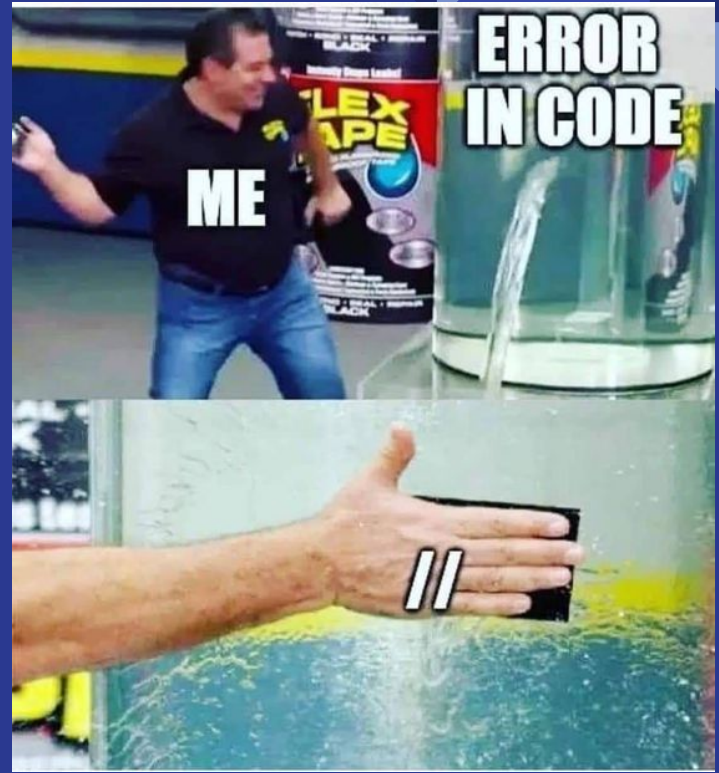
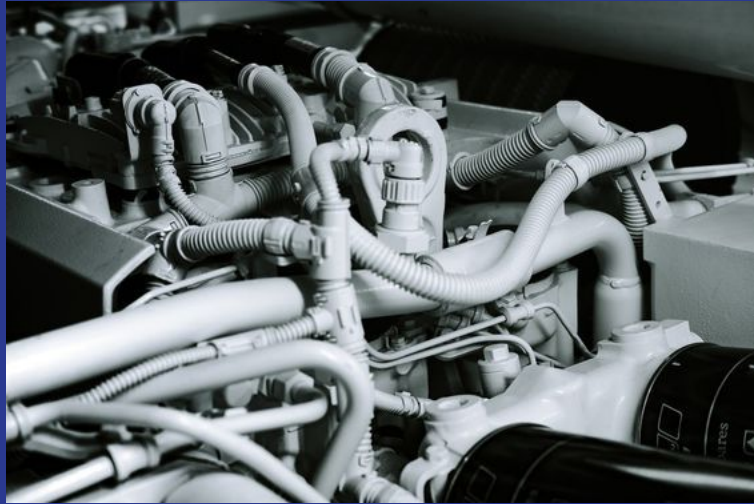
Falácia: A lógica de negócios pode e deve ser centralizada.



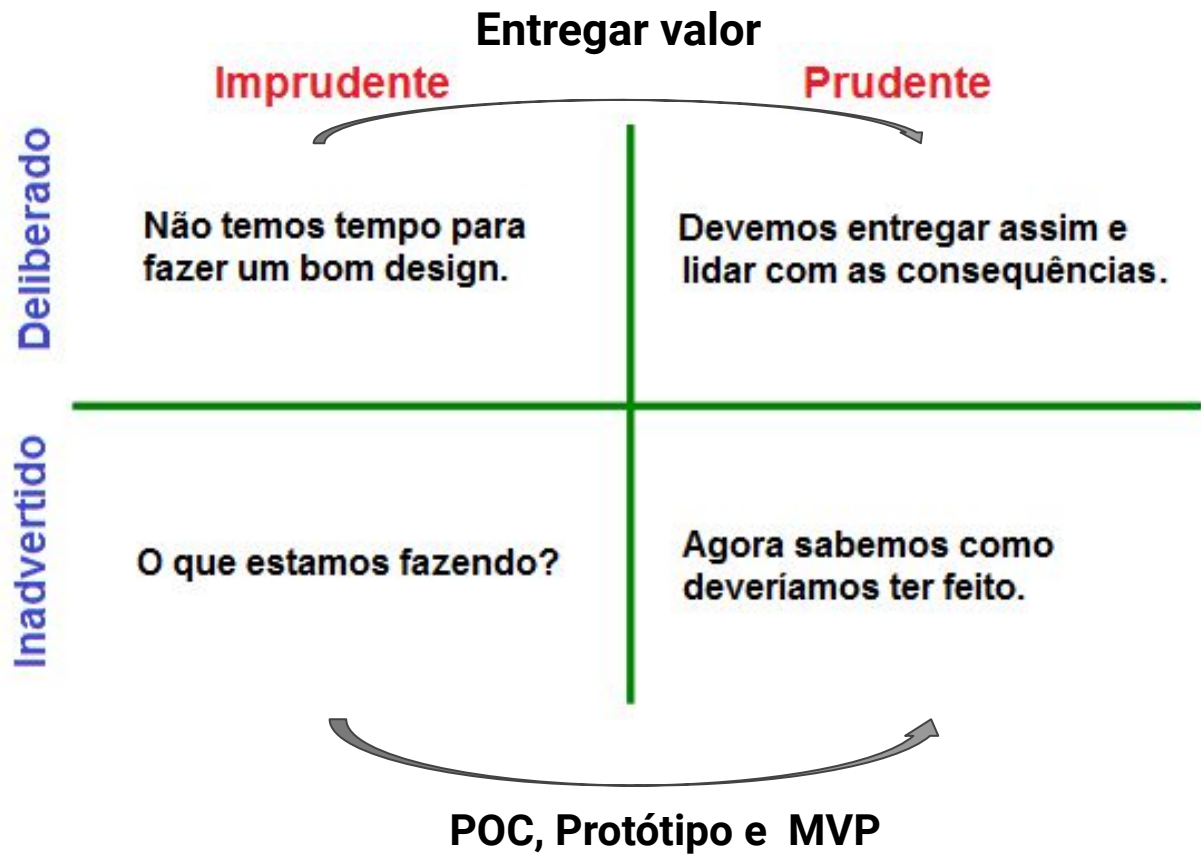


- **Isolar o contexto (DDD - Domain driven design, Data-mesh).**
- **Código limpo ajuda a distribuir o conhecimento no time.**
- **Documentação.**
- **Relacionamento ruim entre áreas.**

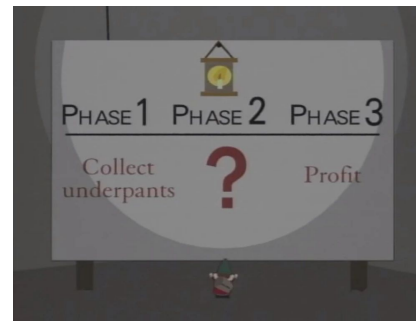
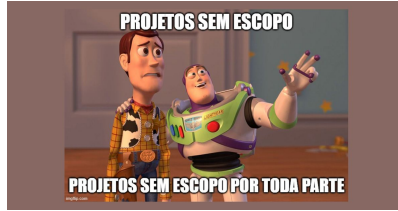
Mitigando a Dívida Técnica



	Imprudente	Prudente
Deliberado	Não temos tempo para fazer um bom design.	Devemos entregar assim e lidar com as consequências.
Inadvertido	O que estamos fazendo?	Agora sabemos como deveríamos ter feito.



Inicie “corretamente”



Desenvolvedor

QUALIDADE	RECURSOS HUMANOS	ESCOPO
AQUISIÇÕES		COMUNICAÇÕES
CUSTO	RISCOS	TEMPO

Definido pelo cliente/negócio/analista (mas não sabe como, quer apenas o resultado)

Já está definido.

Já está definido.

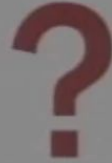


PHASE 1

PHASE 2

PHASE 3

Collect
underpants



Profit

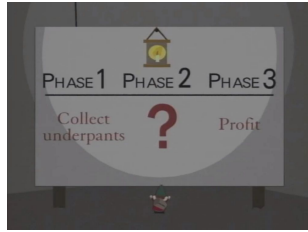




www.dilbert.com
scottadams@aol.com

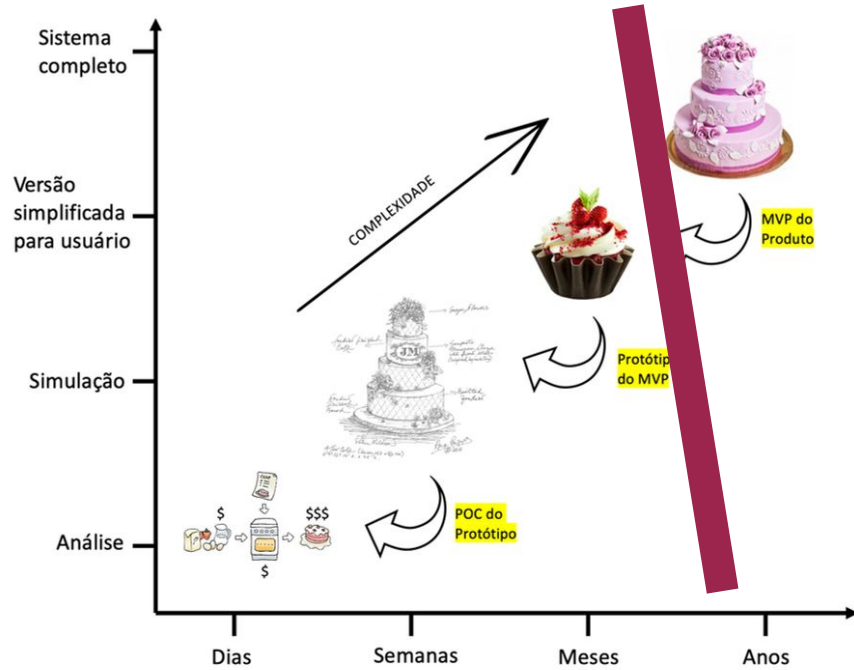
© 2002 United Feature Syndicate, Inc.



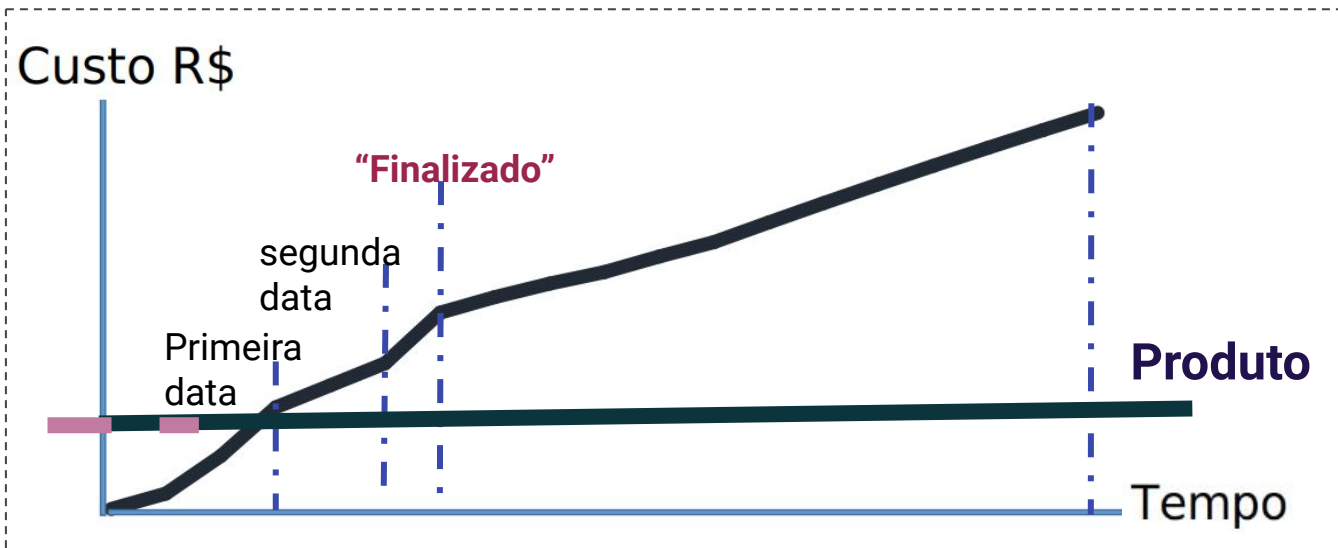


elicitação

Ação ou efeito de elicitar, de fazer sair, de expulsar; elicitação. Obtenção de informações detalhadas sobre o que se pretende fazer.



POC, Protótipo e MVP





Dividir o escopo e eliminar o que não é necessário para entregar valor.

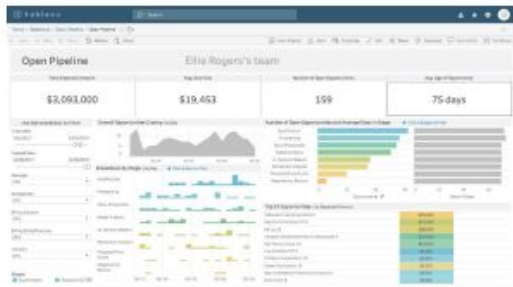
bronze
Table1
Table2
Table3
Table4
Table5
Table6
Table7
Table8
Table9
Table10



silver
Table1
Table2
Table3
Table4
Table5
Table6
Table7
Table8
Table9
Table10



gold
TableX
TableX
TableX
TableX
TableX
TableX
TableX
TableX
TableX
TableX
TableX



bronze
Table1
Table2
Table3
Table4
Table5
Table6
Table7
Table8
Table9
Table10



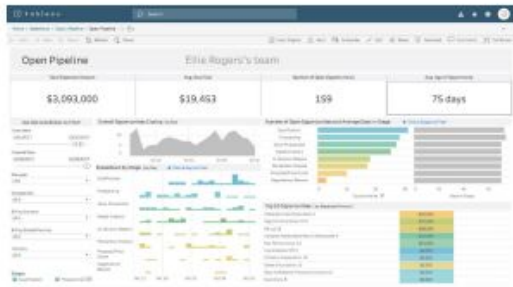
silver
Table1
Table2
Table3
Table4
Table5
Table6
Table7
Table8
Table9
Table10



gold
TableX



gold
TableY

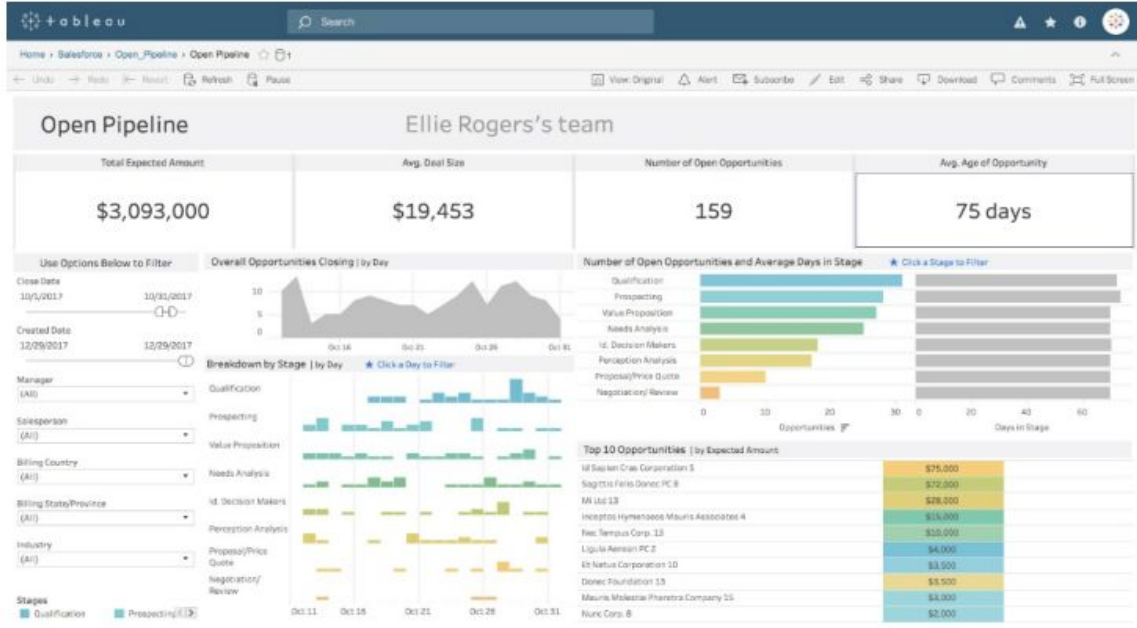


gold

TableY
ColY
ColY
ColY
ColY
ColY
ColY
ColY
ColY
ColY



TableX
ColX
ColX
ColX
ColX
ColX
ColX
ColX
ColX
ColX
ColX

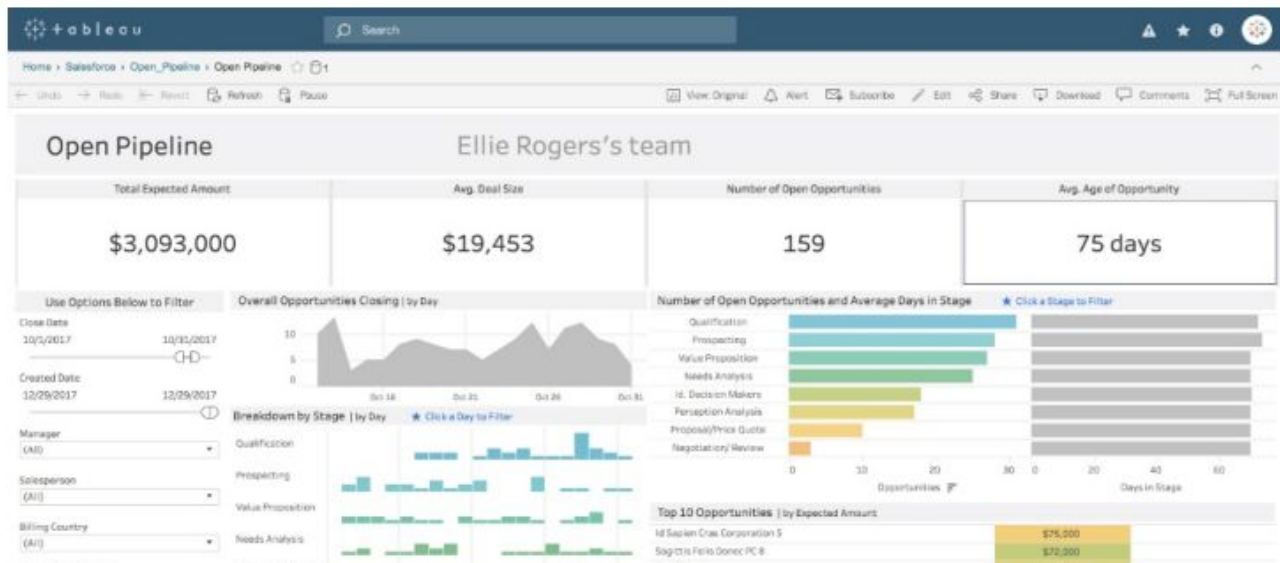


gold

TableY
CoLY
CoLY
CoLY



TableX
CoIX
CoIX

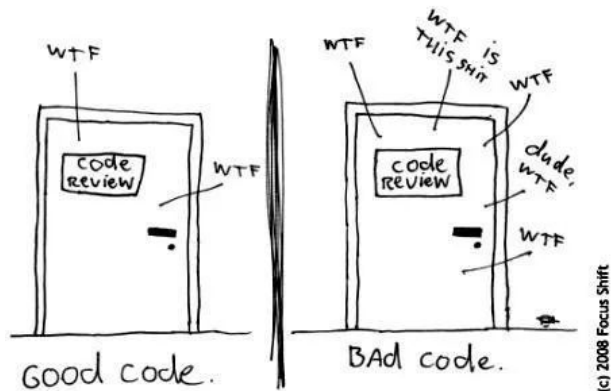


	Tabelas	colunas	colunas (apenas o necessário)
tabelas de sistemas	10	100	
tabelas na gold	2	20	5
	20%	20%	5%
	80% de economia	80% de economia	95% de economia + qualidade

** Profiling (Perfil)

Escreva e **refatore**

The ONLY valid measurement
of code quality: WTFs/minute



Qualquer alteração ou decisão técnica deve ser aplicada como código.

- Infrastructure as Code
- **Data Quality as code**
- Report as code (mais complicado)
- Pipeline as code

Quando refatorar: Feriu a qualidade do dado

Profiling (Perfil):

- máquina de estado (colunas de “status”)
- ids, joins, merge (null, duplicados)
- limites (venda, compra, desconto..)

Single Column

- Cardinalities
- Patterns & Data types
- Value distributions
- Domain classification

Multi-Column

- Correlations
- Association rules
- Clustering
- Outliers
- Summaries & sketches

Cross-DB dependencies



- Unique column combinations
- Inclusion dependencies
- Functional dependencies



Great Expectations

Quando refatorar: código hadouken

```
function register()
{
    if (!empty($_POST)) {
        $msg = '';
        if ($_POST['user_name']) {
            if ($_POST['user_password_new']) {
                if ($_POST['user_password_new'] === $_POST['user_password_repeat']) {
                    if (strlen($_POST['user_password_new']) > 5) {
                        if (strlen($_POST['user_name']) < 65 && strlen($_POST['user_name']) > 1) {
                            if (preg_match('/^[a-z\d]{2,64}$/i', $_POST['user_name'])) {
                                $user = read_user($_POST['user_name']);
                                if (!isset($user['user_name'])) {
                                    if ($_POST['user_email']) {
                                        if (strlen($_POST['user_email']) < 65) {
                                            if (filter_var($_POST['user_email'], FILTER_VALIDATE_EMAIL)) {
                                                create_user();
                                                $_SESSION['msg'] = 'You are now registered so please login';
                                                header('Location: ' . $_SERVER['PHP_SELF']);
                                                exit();
                                            } else $msg = 'You must provide a valid email address';
                                        } else $msg = 'Email must be less than 64 characters';
                                    } else $msg = 'Email cannot be empty';
                                } else $msg = 'Username already exists';
                            } else $msg = 'Username must be only a-z, A-Z, 0-9';
                        } else $msg = 'Username must be between 2 and 64 characters';
                    } else $msg = 'Password must be at least 6 characters';
                } else $msg = 'Passwords do not match';
            } else $msg = 'Empty Password';
        } else $msg = 'Empty Username';
        $_SESSION['msg'] = $msg;
    }
    return register_form();
}
```



Quando refatorar: Código acima de um limite estipulado (50, 100 linhas ?)

```
1 WITH q1 AS (SELECT
2     col1,
3     col2,
4     col3,
5     col4,
6     col5,
7     col6,
8     col7,
9     col8,
10    col9,
11    col10,
12    col11,
13    col12,
14    col13,
15    col14,
16    col15,
17    col16,
18    col17,
19    col18,
20    col19,
21    col20,
22    col21,
23    col22 FROM tablePT01 JOIN tablePT02)
24
25 q1 (SELECT
26     col1,
27     col2,
28     col3,
29     col4,
30     col5,
31     col6,
32     col7,
33     col8,
34     col9,
35     col10,
36     col11,
37     col12,
38     col13,
39     col14,
40     col15,
41     col16,
42     col17,
43     col18,
44     col19,
45     col20,
46     col21,
47     col22 FROM tablePT01 JOIN tablePT02)
48
49 SELECT col1,
50        col2,
51        col3,
52        col4,
53        col5,
54        col6,
55        col7,
56        col8,
57        col9,
58        col10,
59        col11,
60        col12,
61        col13,
62        col14,
63        col15,
64        col16,
65        col17,
66        col18,
67        col19,
68        col20,
69        col21,
70        col22 FROM q1
71
72 UNION ALL
73 SELECT
74     col1,
75     col2,
76     col3,
77     col4,
78     col5,
79     col6,
80     col7,
81     col8,
82     col9,
83     col10,
84     col11,
85     col12,
86     col13,
87     col14,
88     col15,
89     col16,
90     col17,
91     col18,
92     col19,
93     col20,
94     col21,
95     col22 FROM q2 ;
96
```

Inicie “corretamente”:

Usar poc e/ou protótipos para entender e/ou aprender sobre o contexto, descarta assim que possível - **Feedback, deve entregar valor, elicitação.**

“Fatiou passou”:

Levar apenas o necessário (colunas) para a camada mais exposta ao cliente (relatórios/dashboards) – **refatoração/backfill**

Escreva e **refatore**:

Todos os objetos da arquitetura devem estar representados em código, código limpo – **refatore..refatore..refatore.**

Dívida deve ser controlada.

- Dívida deve ser adquirida conscientemente.
- Dívida deve ser **paga**.

Falácias:

- O sistema está finalizado.
- A arquitetura é a mesma em todos os lugares.
- A lógica de negócios pode e deve ser centralizada.

Mitigando a Dívida Técnica:

- **Inicie “corretamente”**
- **“Fatiou passou”**
- **Escreva e refatore**

Dívida deve ser controlada.

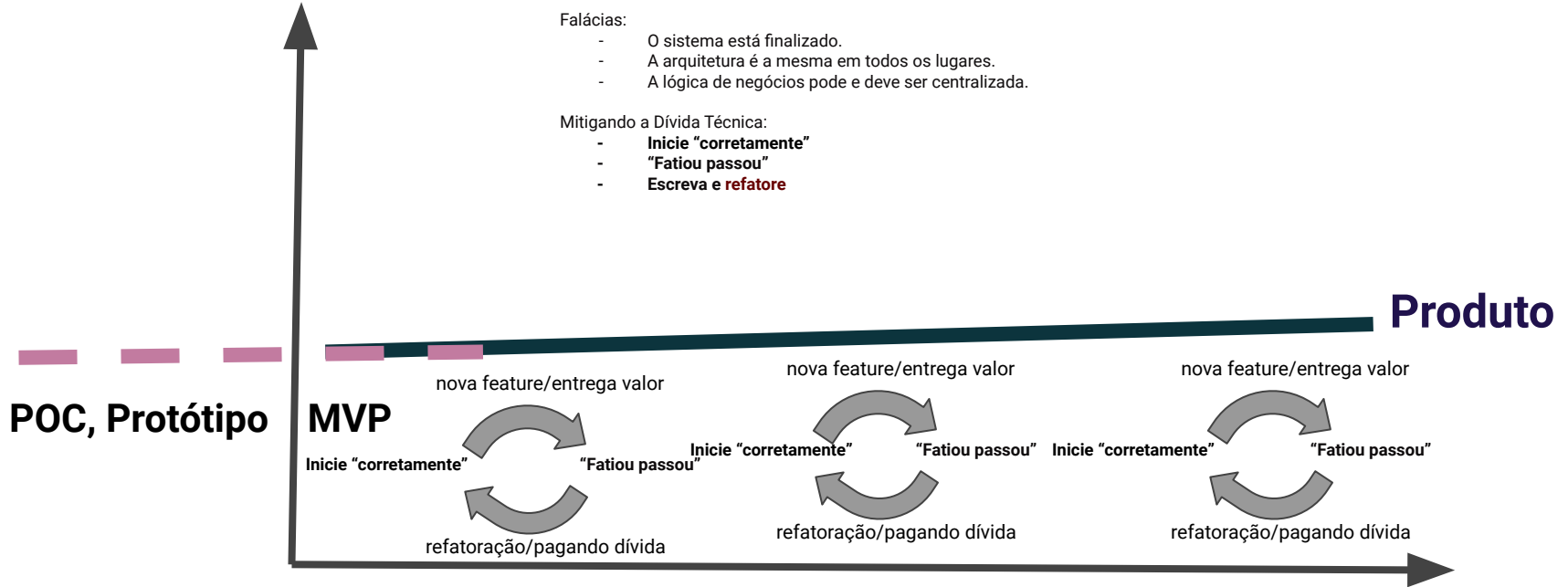
- Dívida deve ser adquirida conscientemente.
- Dívida deve ser **paga**.

Falácias:

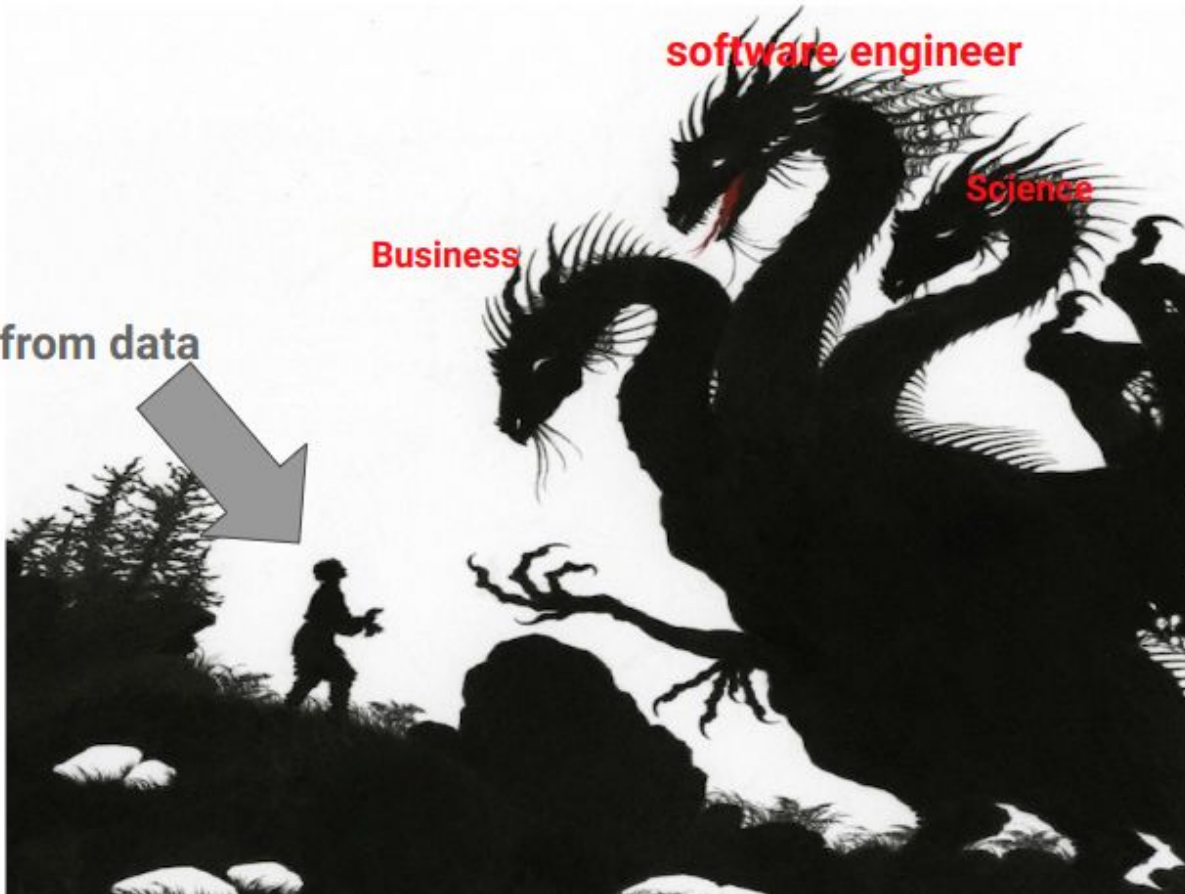
- O sistema está finalizado.
- A arquitetura é a mesma em todos os lugares.
- A lógica de negócios pode e deve ser centralizada.

Mitigando a Dívida Técnica:

- **Inicie "corretamente"**
- **"Fatiou passou"**
- **Escreva e refatore**



someone from data



delivery of
valuable data







Dúvidas ?