# AI Do's and Don'ts

Prolego delivers practical AI engineering services to Fortune 1000 companies

# Getting Real

A firehose of practical advice for your AI journey

# Today's Agenda

- Strategy

- Prototype

- Deployment

Here it comes!

Decide what is relevant for you.

# Strategy

# What can AI do?

# How do we start?



BECOME AN **AI** COMPANY IN 90 DAYS

The complete guide for understanding AI, identifying opportunities, and launching your first product

KEVIN DEWALT

Produced by: Russ Rands

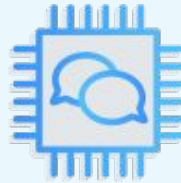**Download an ecopy at book at www.prolego.com**

# AI Product Patterns

**Pattern 1**
Computer Vision

- *What is in this image?*
- *Tell me when something changes*

**Pattern 2**
Natural Language Processing (NLP)

- *Transcribe, translate his speech*
- *Notify analyst if we should review*

**Pattern 3**
Next-in-Sequence Predictions

- *Predicting future sales of an item based on past sales*
- *Predicting which users will buy*
- *Detecting fraudulent credit card transactions*

# Presenting NLP in a New Way: Adventures in AI

*Inspired by GE's use of comic books when introducing new technologies in the 1940s and 1950s, we created our own comic book, Adventures in AI.*

**1950**



**2020**



Download a free copy at
www.prolego.com/adventures-in-ai

# Project Selection

## AI Canvas

### Opportunity
**Why do it?**

A high-level description of the business benefit for the AI models: revenue growth, cost reduction, speed, etc.

### Solution
**What is it?**

A high-level description of the workflow, models, and system architecture.

### Consumers
**Who needs the model outputs?**

AI models produce outputs from input data sources. Consumers are the products, systems, and people who use the outputs to deliver business value.

### Data sources
**What are the model inputs?**

Primary internal/external sources of data for model inputs. Consider complexities such as accessibility, cleansing challenges.

### Strategy
**Why us?**

Unique assets such as customers, data, or business expertise which provide ongoing sustainable advantage.

### Policy & process
**What else must change?**

Necessary data, security, legal, or organizational changes.

### Model development
**How will we build it?**

Identify existing research papers, models, and transfer learning opportunities for accelerating deployment.

### Success criteria
**How will we know it works?**

KPIs or other business metrics to gauge success. Qualitative feedback as necessary.

# Prototype

# Get Real!

**Budget:** *creating capability, not immediate ROI*

**Team:** *current talent before recruiting*

**Data:** *don't delay due to "data issues"*

# Evaluate impact on your data

# Tips

- Time-box to 4 months max!

- Work with engaged business partners

- Plan for deployment on day 0

# Deployment

That light
… is indeed
a train

# From Notebooks to Production Code

# Closing Thoughts

- Think defense.

- Think long term.

- Think capabilities - not solutions.