

Generative Interpretable Visual Design

Application to Visual Conjoint

Ankit Sisodia¹, Alex Burnap² and Vineet Kumar²

¹Purdue University

²Yale School of Management

Presenting at: ORSI ICBAI 2023

Visual (or aesthetic) design matters across many product categories . . .



Cars



Fashion



Furniture

Visual design matters



Visual design matters



“Exterior look/design is the top reason shoppers avoid a particular vehicle (30%), followed by cost (17%).”

—JD Power Avoider Study 2015

What this paper seeks to do

Research Goals

Our research aims to obtain **interpretable** visual characteristics (not surprising / outlier) directly from unstructured product images

- *automatically discover (extract) characteristics*

What this paper seeks to do

Research Goals

Our research aims to obtain **interpretable** visual characteristics (not surprising / outlier) directly from unstructured product images

- *automatically discover* (extract) characteristics
- *quantify* these characteristics

What this paper seeks to do

Research Goals

Our research aims to obtain **interpretable** visual characteristics (not surprising / outlier) directly from unstructured product images

- *automatically discover* (extract) characteristics
- *quantify* these characteristics
- *generate visual design that span the space of visual characteristics*

What this paper seeks to do

Research Goals

Our research aims to obtain **interpretable** visual characteristics (not surprising / outlier) directly from unstructured product images

- *automatically discover* (extract) characteristics
- *quantify* these characteristics
- *generate* visual design that span the space of visual characteristics

What this paper seeks to do

Research Goals

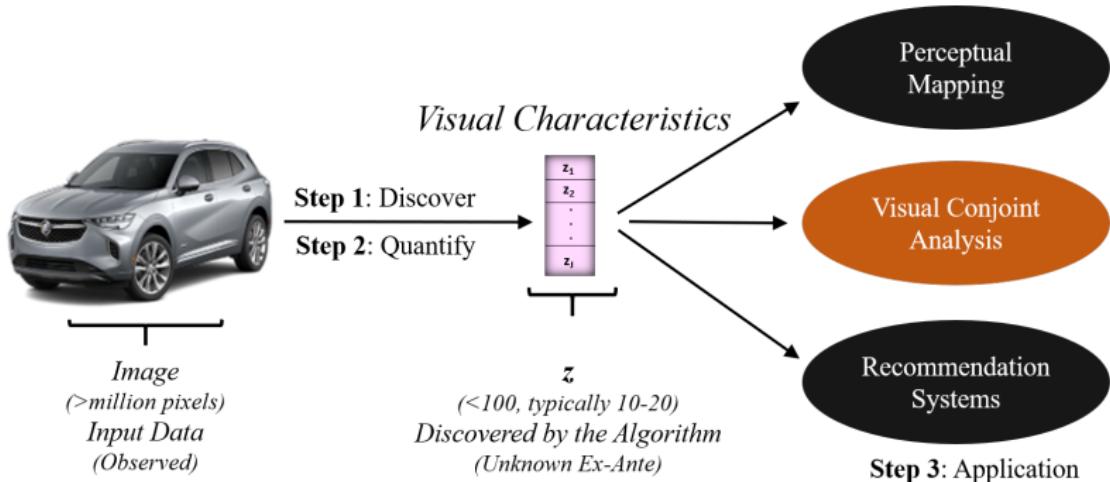
Our research aims to obtain **interpretable** visual characteristics (not surprising / outlier) directly from unstructured product images

- *automatically discover* (extract) characteristics
- *quantify* these characteristics
- *generate* visual design that span the space of visual characteristics

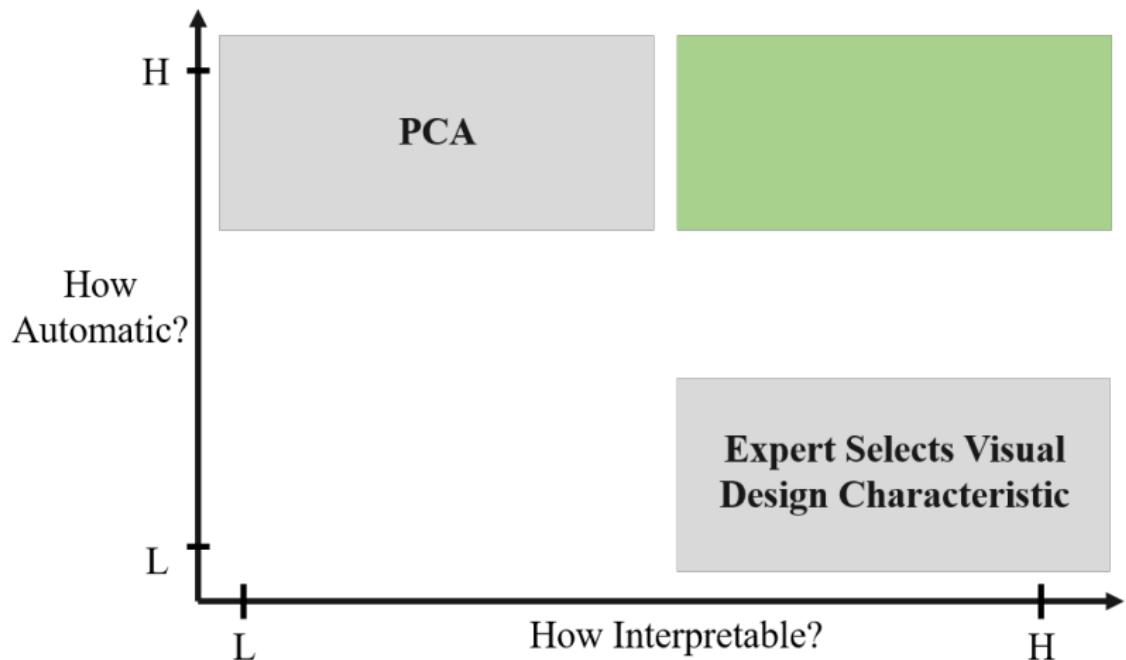


Hyundai: (3, 8, 5, 9) compared to BMW: (1, 3, 10, 1)

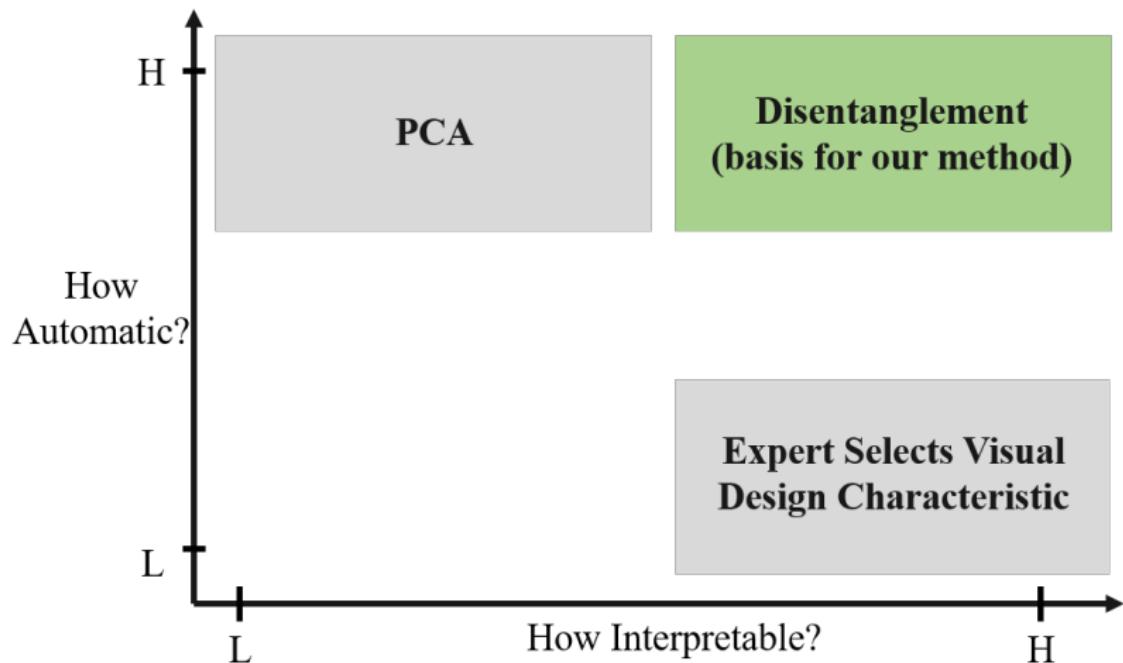
Why Visual Characteristics?



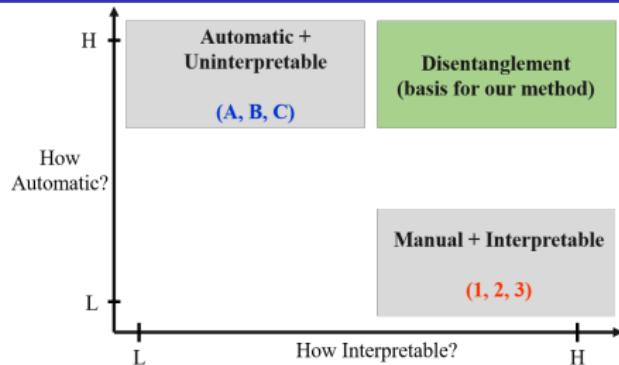
Modeling Visual Characteristics: A comparison of methods



Modeling Visual Characteristics: A comparison of methods



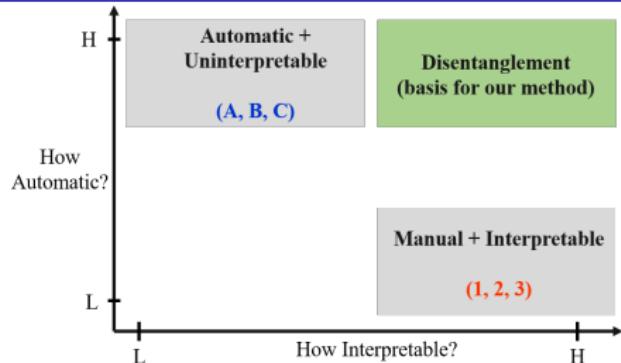
Modeling Visual Characteristics: A comparison of methods



Automatic + Uninterpretable

- A - Bajari, P. L. et al. (2021) : Hedonic prices and quality adjusted price indices powered by AI, *CENMAP working paper*
- B - Law, S., et al. (2019) : Take a look around: using street view and satellite images to estimate house prices. *ACM Transactions on Intelligent Systems and Technology (TIST)*
- C - Aubry, S., et al. (2019) : Machine learning, human experts, and the valuation of real assets. *CFS Working Paper Series*

Modeling Visual Characteristics: A comparison of methods



Automatic + Uninterpretable

- A - Bajari, P. L. et al. (2021) : Hedonic prices and quality adjusted price indices powered by AI, *CENMAP working paper*
- B - Law, S., et al. (2019) : Take a look around: using street view and satellite images to estimate house prices. *ACM Transactions on Intelligent Systems and Technology (TIST)*
- C - Aubry, S., et al. (2019) : Machine learning, human experts, and the valuation of real assets. *CFS Working Paper Series*

Manual + Interpretable

- 1 - Zhang, M. et al. (2022) : Can consumer-posted photos serve as a leading indicator of restaurant survival? Evidence from yelp. *Management Science*
- 2 - Liu, Y., et al. (2017) : The effects of products' aesthetic design on demand and marketing-mix effectiveness: The role of segment prototypicality and brand consistency. *Journal of Marketing*
- 3 - Zhang, S., et al. (2021) : What makes a good image? Airbnb demand analytics leveraging interpretable image features. *Management Science*

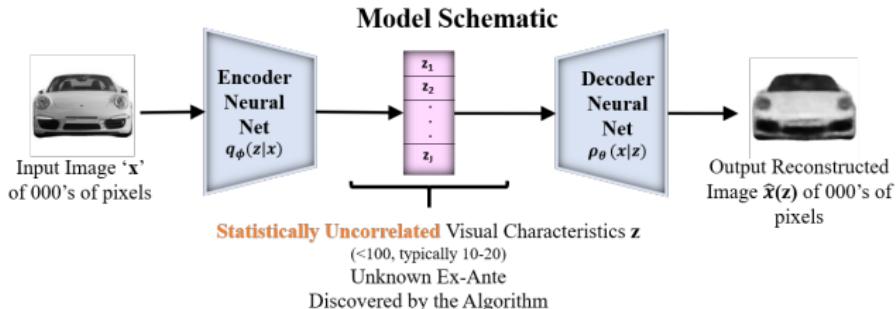
What is disentanglement?

Bengio et al (2013)

*"A disentangled representation can be defined as one where **single latent units** are sensitive to changes in **single generative factors**, while being relatively invariant to changes in other factors"*

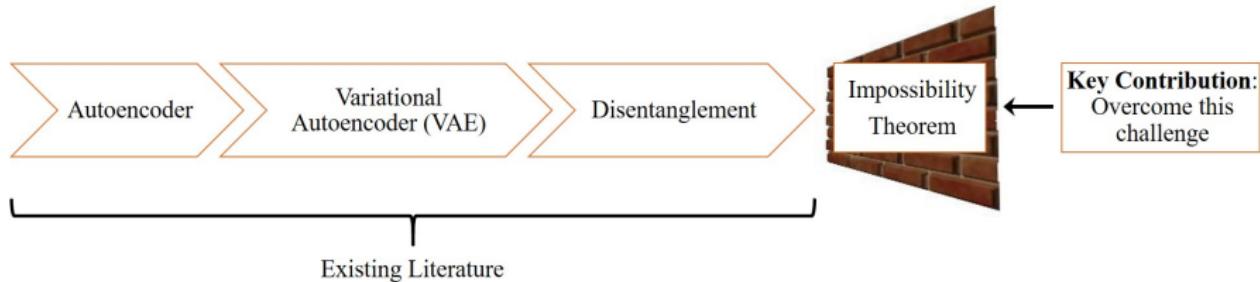
- Latent Units (\mathbf{z}): Dimensions in the model's latent space
- Generative factors (\mathbf{c}): Human-interpretable true characteristics

Models in Existing Literature



Model	Goal
Autoencoder (AE)	Reconstruction accuracy
Variational Autoencoder (VAE)	...+ structured latent space
Disentanglement	...+ ...+ statistically independent latent space

Roadmap of Our Approach



Contribution

We aim to overcome this impossibility theorem by using structured product characteristics.

Disentangled and Entangled Representations

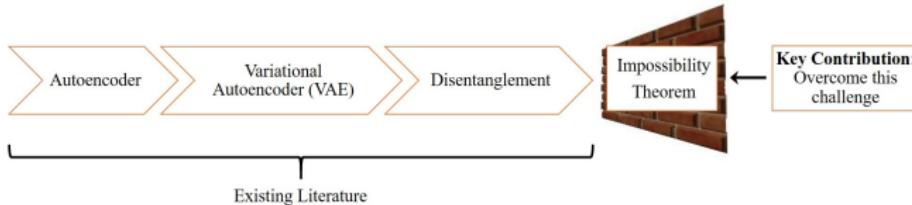
Example of *Entangled* Visual Characteristics



Example of *Disentangled* Visual Characteristics



Impossibility Theorem



Impossibility Theorem

Unsupervised (*i.e. only images*) learning of disentangled representations is *fundamentally impossible* except under certain restrictive conditions.^a

^aLocatello, Francesco, et al. "Challenging common assumptions in the unsupervised learning of disentangled representations." ICML. PMLR, 2019.

Implication: Every disentangled representation can have other *infinite* equivalent entangled representations.

Impossibility Theorem – Implications



z_1
z_2
.
.
.
z_j

predicts →

A horizontal arrow pointing from the learned characteristics table to the ground truth characteristics table.

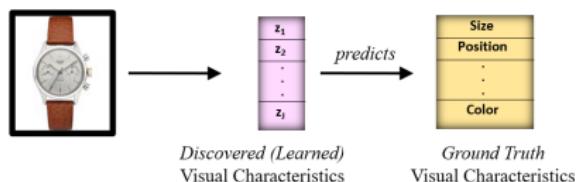
Size
Position
.
.
.
Color

Discovered (Learned)
Visual Characteristics

Ground Truth
Visual Characteristics

Impossibility Theorem – Implications

Common approach to ground truth in ML is to get humans to label¹

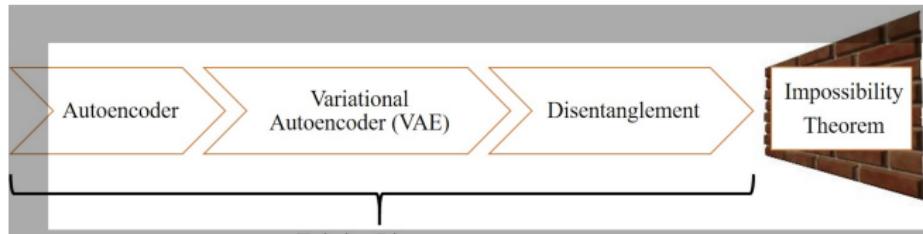


What's the Problem?

- Ground truth on visual characteristics is *unknown*. In fact, these are precisely what we want to find.

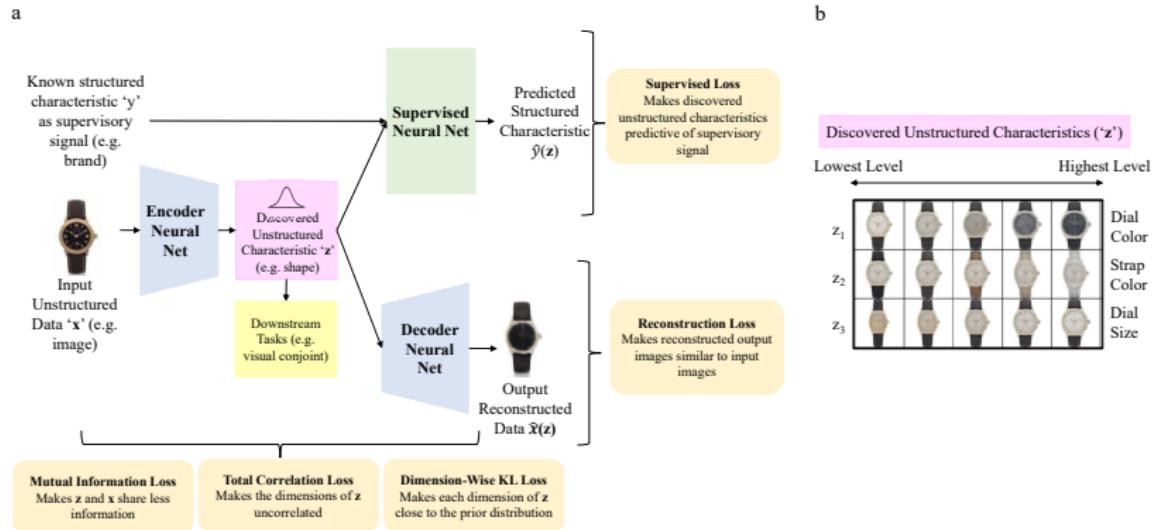
¹Locatello, Francesco, et al. "Disentangling factors of variation using few labels." ICLR. 2020.

Contribution



- **Solution** without ground truth on visual characteristics:
- Leverage **structured product characteristics** to provide a supervisory signal for disentanglement

Schematic of Proposed Approach



Model

- Learn model parameters by minimizing loss $L(\theta, \phi; \mathbf{x}, \mathbf{z})$ of integrated model
- θ and ϕ are encoder and decoder parameters; \mathbf{x} are images

$$\underbrace{L(\theta, \phi, \mathbf{w}; \mathbf{x}, \mathbf{z})}_{\text{Total Loss}} = \underbrace{\mathbf{E}_{q_\phi(\mathbf{z}|\mathbf{x})} [\log p_\theta(\mathbf{x}|\mathbf{z})]}_{\text{Reconstruction Loss}} + \alpha \underbrace{I_q(\mathbf{z}, \mathbf{x})}_{\text{Mutual Information Loss}} + \beta \underbrace{KL \left[q(\mathbf{z}) || \prod_{j=1}^J q(z_j) \right]}_{\text{Total Correlation Loss}} \\ + \gamma \underbrace{\sum_{j=1}^J KL \left[q(z_j) || p(z_j) \right]}_{\text{Dimension-Wise KL Divergence Loss}} + \delta \underbrace{P(\hat{\mathbf{y}}(\mathbf{z}), \mathbf{y})}_{\text{Supervised Loss}}$$

Model

- Learn model parameters by minimizing loss $L(\theta, \phi; \mathbf{x}, \mathbf{z})$ of integrated model
- θ and ϕ are encoder and decoder parameters; \mathbf{x} are images

$$\underbrace{L(\theta, \phi, \mathbf{w}; \mathbf{x}, \mathbf{z})}_{\text{Total Loss}} = \underbrace{\mathbf{E}_{q_\phi(\mathbf{z}|\mathbf{x})} [\log p_\theta(\mathbf{x}|\mathbf{z})]}_{\text{Reconstruction Loss}} + \alpha \underbrace{I_q(\mathbf{z}, \mathbf{x})}_{\text{Mutual Information Loss}} + \beta \underbrace{KL \left[q(\mathbf{z}) || \prod_{j=1}^J q(z_j) \right]}_{\text{Total Correlation Loss}} \\ + \gamma \underbrace{\sum_{j=1}^J KL \left[q(z_j) || p(z_j) \right]}_{\text{Dimension-Wise KL Divergence Loss}} + \delta \underbrace{P(\hat{\mathbf{y}}(\mathbf{z}), \mathbf{y})}_{\text{Supervised Loss}}$$

Loss Term	Why is this term included?
Reconstruction	Promotes accurate reconstruction of images
Mutual Information	Minimizes redundant information
Total Correlation	Promotes statistical independence between visual characteristics
Dimension-Wise KL	Penalizes deviations from a prior
Supervised	Provides a signal to address the impossibility theorem

Model – Role of Supervised Loss

$$\underbrace{L(\theta, \phi, \mathbf{w}; \mathbf{x}, \mathbf{z})}_{\text{Total Loss}} = \underbrace{\mathbf{E}_{q_\phi(\mathbf{z}|\mathbf{x})} [\log p_\theta(\mathbf{x}|\mathbf{z})]}_{\text{Reconstruction Loss}} + \alpha \underbrace{I_q(\mathbf{z}, \mathbf{x})}_{\text{Mutual Information Loss}} + \beta \underbrace{KL \left[q(\mathbf{z}) || \prod_{j=1}^J q(z_j) \right]}_{\text{Total Correlation Loss}} \\ + \gamma \underbrace{\sum_{j=1}^J KL \left[q(z_j) || p(z_j) \right]}_{\text{Dimension-Wise KL Divergence Loss}} + \delta \underbrace{P(\hat{\mathbf{y}}(\mathbf{z}), \mathbf{y})}_{\text{Supervised Loss}}$$

- Supervised Loss is used to predict signal from latent representation z : $s = f(z)$

Model – Role of Supervised Loss

$$\underbrace{L(\theta, \phi, \mathbf{w}; \mathbf{x}, \mathbf{z})}_{\text{Total Loss}} = \underbrace{\mathbf{E}_{q_\phi(\mathbf{z}|\mathbf{x})} [\log p_\theta(\mathbf{x}|\mathbf{z})]}_{\text{Reconstruction Loss}} + \alpha \underbrace{I_q(\mathbf{z}, \mathbf{x})}_{\text{Mutual Information Loss}} + \beta \underbrace{KL \left[q(\mathbf{z}) || \prod_{j=1}^J q(z_j) \right]}_{\text{Total Correlation Loss}} \\ + \gamma \underbrace{\sum_{j=1}^J KL \left[q(z_j) || p(z_j) \right]}_{\text{Dimension-Wise KL Divergence Loss}} + \delta \underbrace{P(\hat{\mathbf{y}}(\mathbf{z}), \mathbf{y})}_{\text{Supervised Loss}}$$

- Supervised Loss is used to predict signal from latent representation z : $s = f(z)$
- Can use structured product characteristics as signals: brand, price, material etc.

Model – Role of Supervised Loss

$$\underbrace{L(\theta, \phi, \mathbf{w}; \mathbf{x}, \mathbf{z})}_{\text{Total Loss}} = \underbrace{\mathbf{E}_{q_\phi(\mathbf{z}|\mathbf{x})} [\log p_\theta(\mathbf{x}|\mathbf{z})]}_{\text{Reconstruction Loss}} + \alpha \underbrace{I_q(\mathbf{z}, \mathbf{x})}_{\text{Mutual Information Loss}} + \beta \underbrace{KL \left[q(\mathbf{z}) || \prod_{j=1}^J q(z_j) \right]}_{\text{Total Correlation Loss}} \\ + \gamma \underbrace{\sum_{j=1}^J KL \left[q(z_j) || p(z_j) \right]}_{\text{Dimension-Wise KL Divergence Loss}} + \delta \underbrace{P(\hat{\mathbf{y}}(\mathbf{z}), \mathbf{y})}_{\text{Supervised Loss}}$$

- Supervised Loss is used to predict signal from latent representation z : $s = f(z)$
- Can use structured product characteristics as signals: brand, price, material etc.

Model – Role of Supervised Loss

$$\underbrace{L(\theta, \phi, \mathbf{w}; \mathbf{x}, \mathbf{z})}_{\text{Total Loss}} = \underbrace{\mathbf{E}_{q_\phi(\mathbf{z}|\mathbf{x})} [\log p_\theta(\mathbf{x}|\mathbf{z})]}_{\text{Reconstruction Loss}} + \alpha \underbrace{I_q(\mathbf{z}, \mathbf{x})}_{\text{Mutual Information Loss}} + \beta \underbrace{KL \left[q(\mathbf{z}) || \prod_{j=1}^J q(z_j) \right]}_{\text{Total Correlation Loss}} \\ + \gamma \underbrace{\sum_{j=1}^J KL \left[q(z_j) || p(z_j) \right]}_{\text{Dimension-Wise KL Divergence Loss}} + \delta \underbrace{P(\hat{\mathbf{y}}(\mathbf{z}), \mathbf{y})}_{\text{Supervised Loss}}$$

- Supervised Loss is used to predict signal from latent representation z : $s = f(z)$
- Can use structured product characteristics as signals: brand, price, material etc.

Idea to Overcome Impossibility Theorem

If the supervisory signal is sufficiently correlated with visual characteristics, then it can help obtain the unique (true) disentangled representation

Human Interpretable Characteristics?

- UDR indicates disentanglement, but are these visual characteristics human interpretable?
 - Without any domain knowledge about the product category?

Human Interpretable Characteristics?

- UDR indicates disentanglement, but are these visual characteristics human interpretable?
 - Without any domain knowledge about the product category?



Starting from the image on the left, what part of the watch changes the most as you go from left to right? Carefully check both large and small visual aspects. Go through each part of the watch one by one before selecting any option. Refer to the above image to see parts of the watch.



Note: Images are low-quality on purpose

- | | |
|-----------------------------------|-----------------------------------|
| <input type="radio"/> Bezel | <input type="radio"/> Hands |
| <input type="radio"/> Crown | <input type="radio"/> Hour Marker |
| <input type="radio"/> Date Window | <input type="radio"/> Lug |
| <input type="radio"/> Dial | <input type="radio"/> Strap |

How is that part of the watch changing?

Visual Characteristics: Quantification?

Interpretability and Quantification

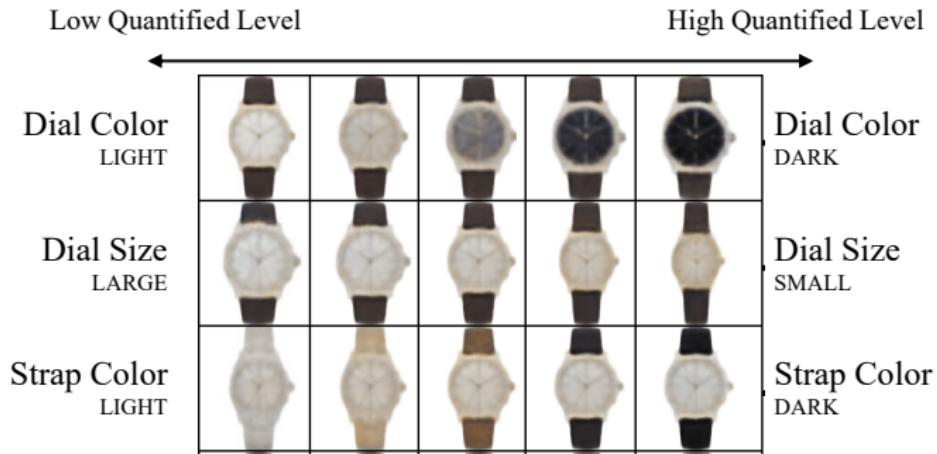
Visual characteristic	Interpretability Survey	Quantification Survey
Dial Size	76%	83%
Dial Color	80%	92%
Strap Color	88%	92%
Rim (Bezel) Color	79%	88%
Dial Shape	87%	68%
Knob (Crown) Size	70%	85%

Discovered Visual characteristics

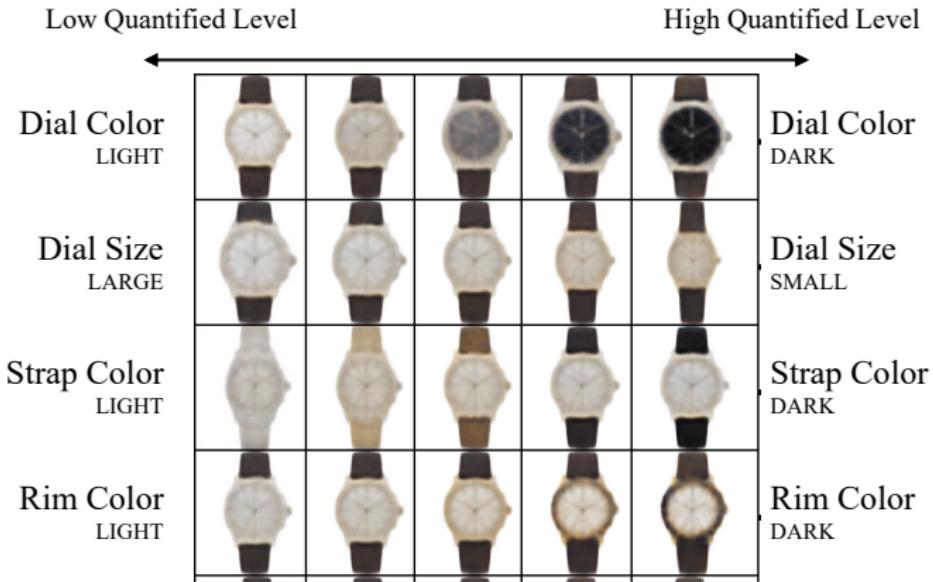
Discovered Visual characteristics



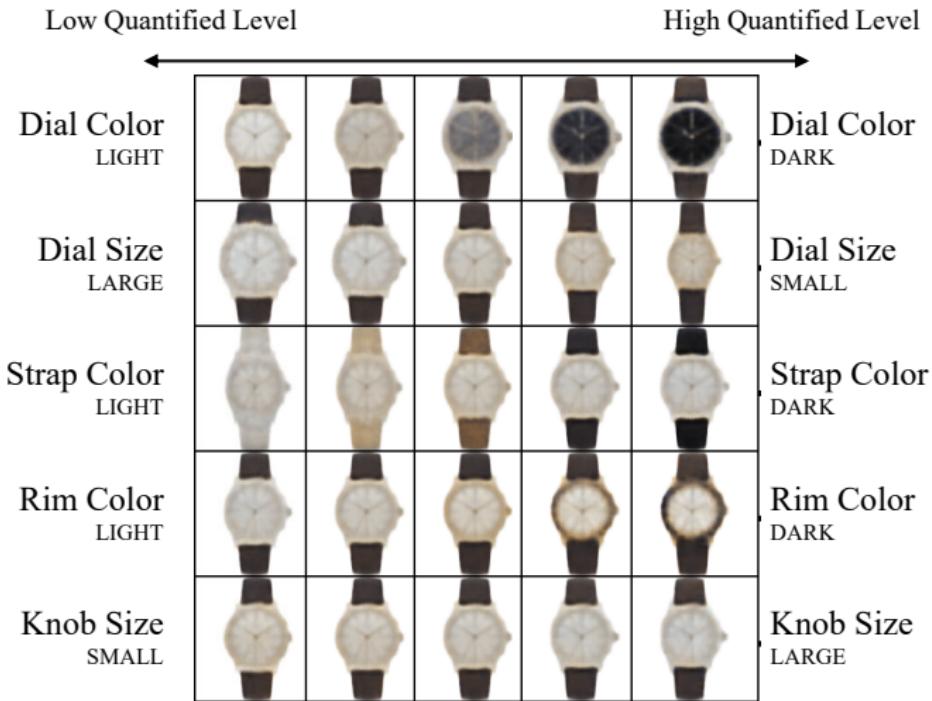
Discovered Visual characteristics



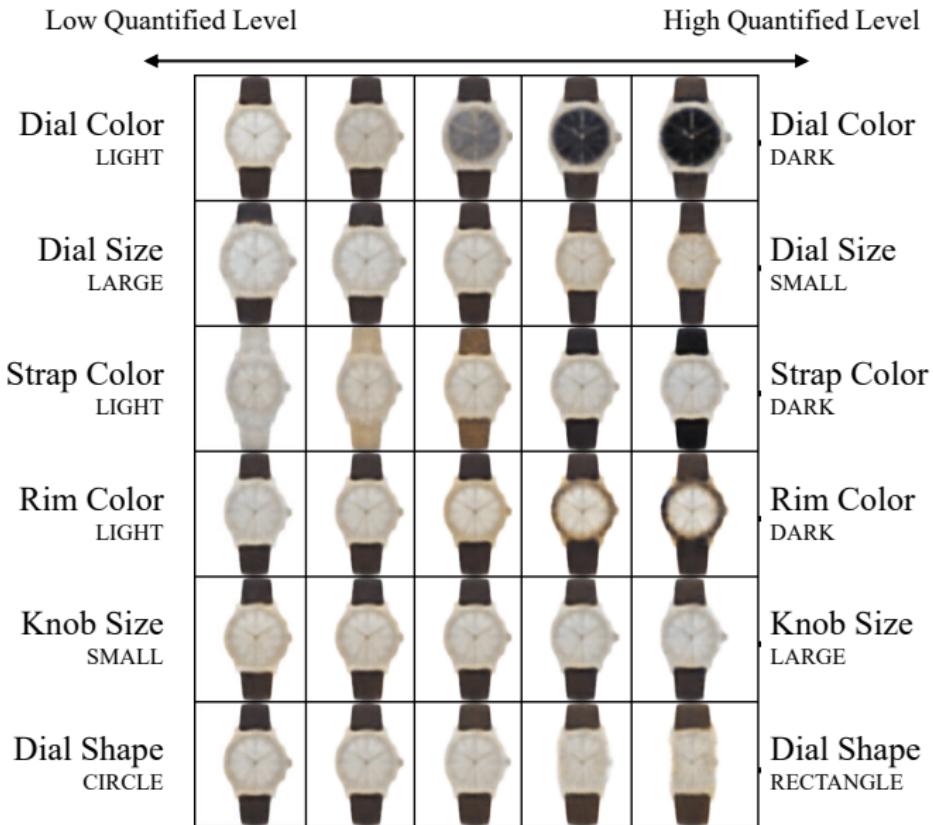
Discovered Visual characteristics



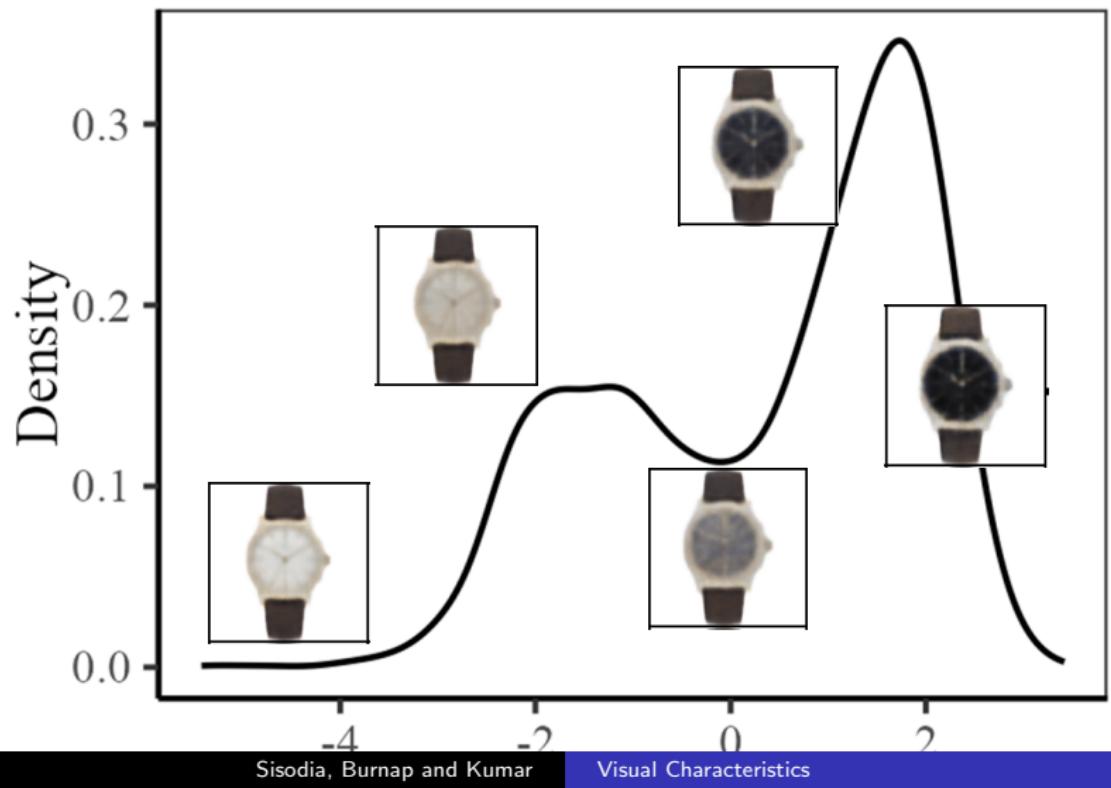
Discovered Visual characteristics



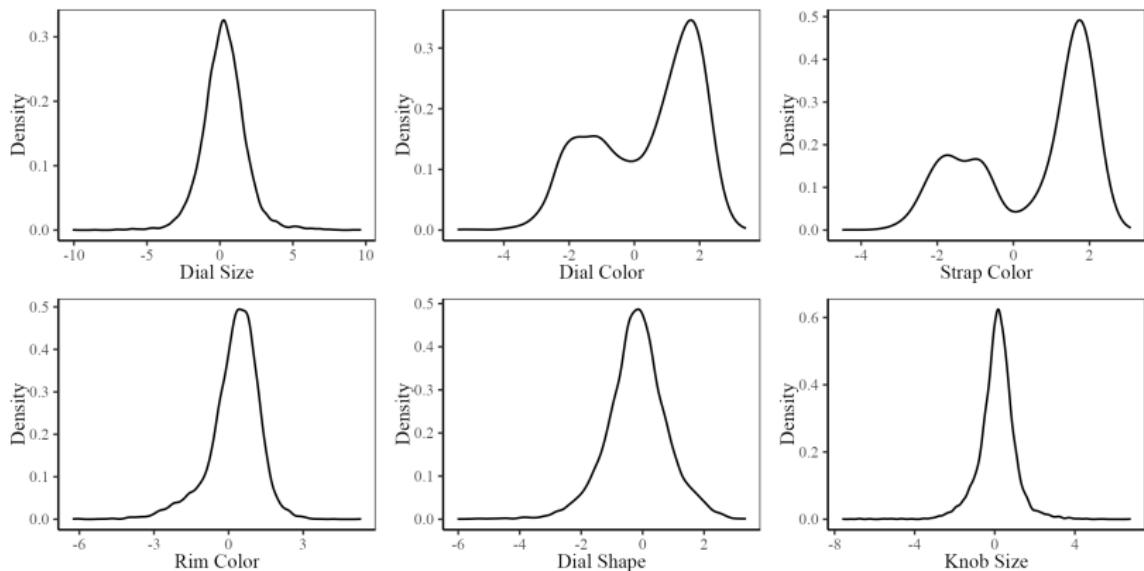
Discovered Visual characteristics



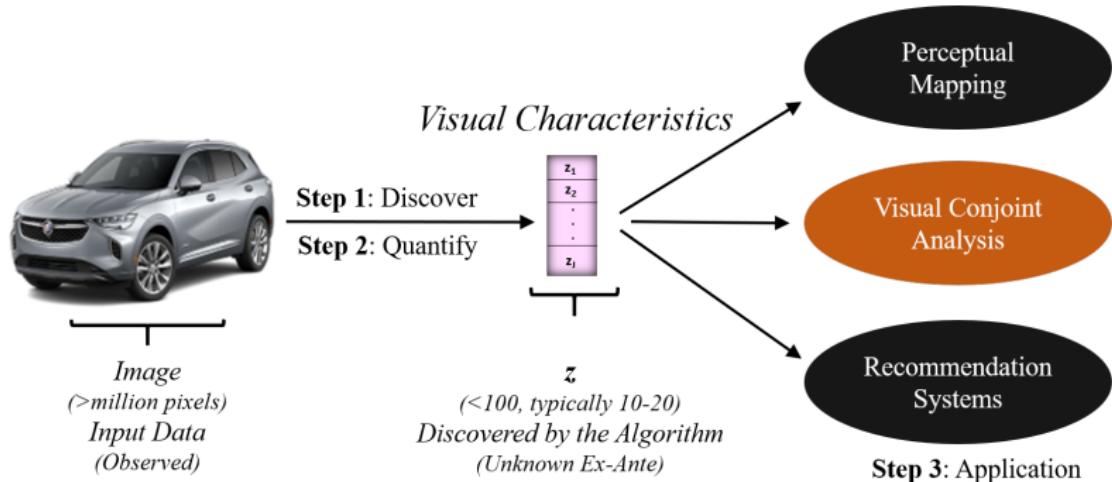
Density of Discovered Visual characteristics (from 'Brand+Material' Signal)



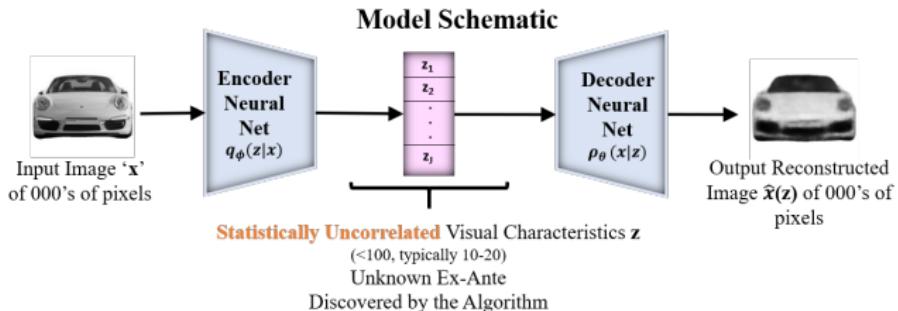
Density of Discovered Visual characteristics (from 'Brand+Material' Signal)



Research Goals

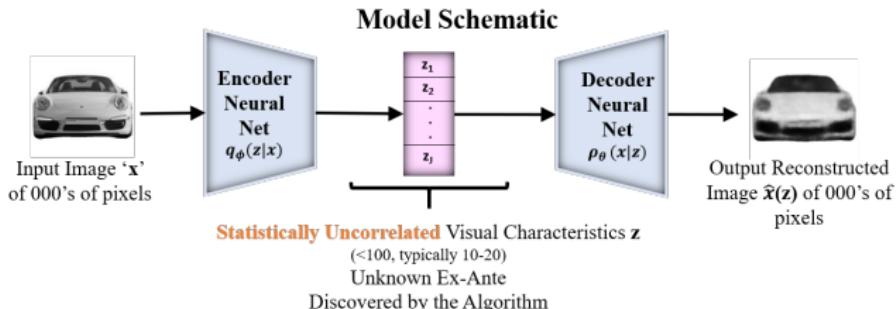


Visual Conjoint Analysis: Background



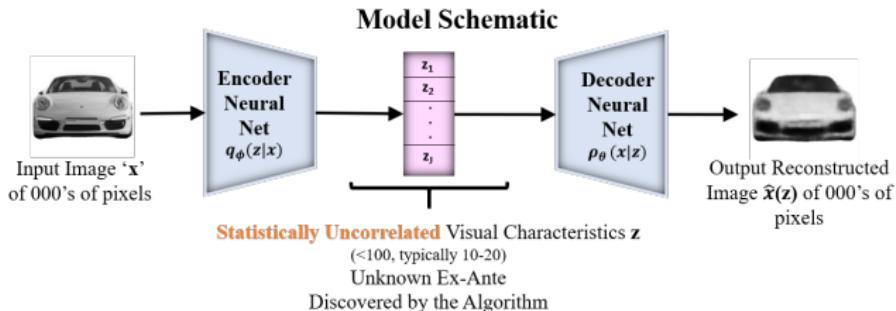
- Visual conjoint has been challenging to do because elements of visual space are correlated

Visual Conjoint Analysis: Background



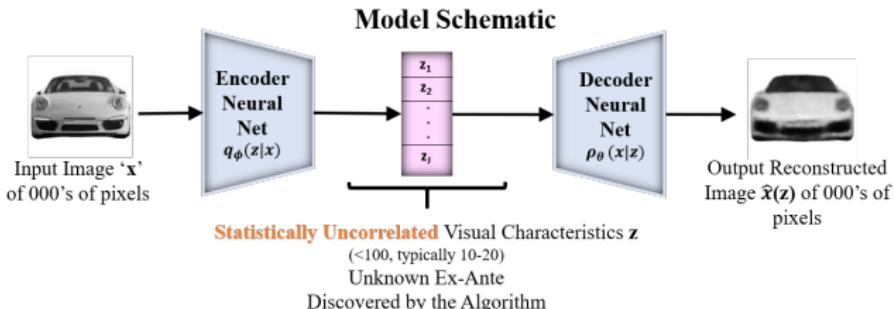
- Visual conjoint has been challenging to do because elements of visual space are correlated
- Designs have always been manually generated by product designers (prototypes)

Visual Conjoint Analysis: Background



- Visual conjoint has been challenging to do because elements of visual space are correlated
- Designs have always been manually generated by product designers (prototypes)
- Our approach generates new never-seen visual designs (counterfactual)

Visual Conjoint Analysis: Background



- Visual conjoint has been challenging to do because elements of visual space are correlated
- Designs have always been manually generated by product designers (prototypes)
- Our approach generates new never-seen visual designs (counterfactual)
- **Can span the entire space of visual designs *without being bound by the correlations in the data.***

Conclusion

We obtain interpretable visual characteristics directly from unstructured product images

- *automatically discover (extract) characteristics*

Applications

We then used the model to:

- generate new counterfactual designs to obtain consumer preferences over visual characteristics.
- obtain ideal point visual designs corresponding to different consumer segments

Conclusion

We obtain interpretable visual characteristics directly from unstructured product images

- *automatically discover* (extract) characteristics
- *quantify these characteristics*

Applications

We then used the model to:

- generate new counterfactual designs to obtain consumer preferences over visual characteristics.
- obtain ideal point visual designs corresponding to different consumer segments

Conclusion

We obtain interpretable visual characteristics directly from unstructured product images

- automatically discover (extract) characteristics
- quantify these characteristics
- generate visual design that span the space of visual characteristics

Applications

We then used the model to:

- generate new counterfactual designs to obtain consumer preferences over visual characteristics.
- obtain ideal point visual designs corresponding to different consumer segments