

Rent3D:

Floor-Plan Priors for Monocular Layout Estimation

Chenxi Liu^{1,*} Alexander Schwing^{2,*} Kaustav Kundu²
Raquel Urtasun² Sanja Fidler²

¹Tsinghua University, ²University of Toronto



How Many Times Have You Looked for Apartments?



How Many Times Have You Looked for Apartments?



United States:

- 11.7% per year

Craigslist:

- 90,000 rental ads per day only in New York
- 10 million people visit the website per day

How Many Times Have You Looked for Apartments?



2 times



3 times



2 times



4 times



5 times

Finding an Apartment/House is a Pain...

- Particularly during a winter in Toronto



Renting Apartments

5 bedroom apartment for sale

£64,999,950

One Hyde Park, Knightsbridge, SW1X



Start slideshow

1 of 10

Enlarge Picture No.02



Do you like this property?

Call: 020 8012 4022

Request Details

Description

Floorplan

Map & Schools

Street View

Virtual Tour

Full description

« Back to property listings

This property is marketed by:



Aylesford International, Chelsea
440 Kings Road, London, SW10 0LH

View properties from this agent

Request Details

or call: 020 8012 4022

★ Save property

📝 Add notes

🖨️ Print

✉️ Send to Friend

Share this property

f Share

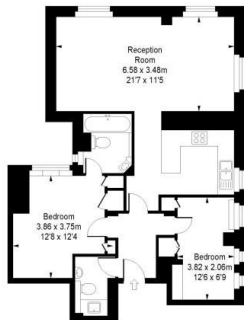
🐦 Tweet

Pin it

don't miss out

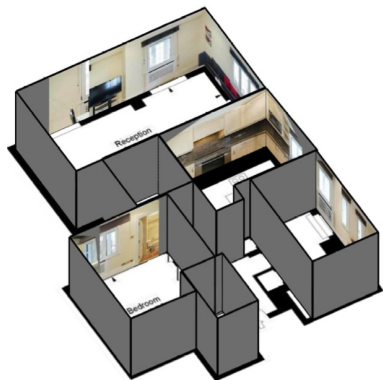
75% of home-movers in

Example Rental Data

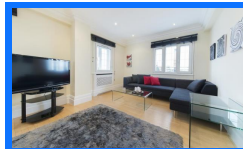
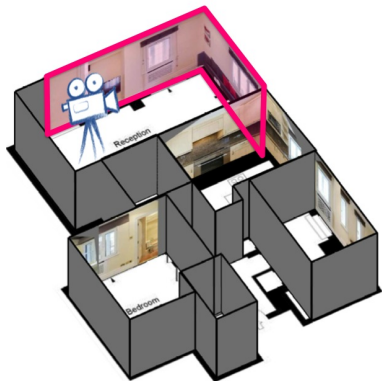


- Plus some meta information e.g. wall height

Rent3D: View Rental Ads in 3D



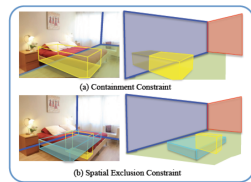
Rent3D: View Rental Ads in 3D



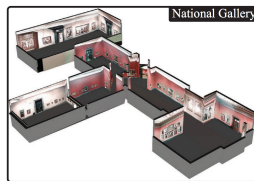
- Camera localization within apartment

Related Work

- Room layout estimation
 - ▷ Hedau et al., 2009, 2012
 - ▷ Lee et al., 2010
 - ▷ Schwing et al., 2012, 2013
 - ▷ Del Pero et al., 2011, 2012
 - ▷ Choi et al., 2013
- Virtual tours
 - ▷ Xiao & Furukawa, 2012
- 3D indoor reconstruction from large photo collections or video
 - ▷ Cabral & Furukawa, 2014
 - ▷ Brualla et al., 2014
- Indoor localization (video, depth sensors)
 - Project Tango
 - SLAM work



Lee et al., 2010

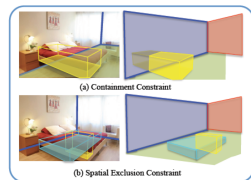


Xiao & Furukawa, 2012



Cabral & Furukawa, 2014

- Room layout estimation
 - ▷ Hedau et al., 2009, 2012
 - ▷ Lee et al., 2010
 - ▷ Schwing et al., 2012, 2013
 - ▷ Del Pero et al., 2011, 2012
 - ▷ Choi et al., 2013



Lee et al., 2010

Our work: 3D indoor reconstruction and localization using monocular imagery

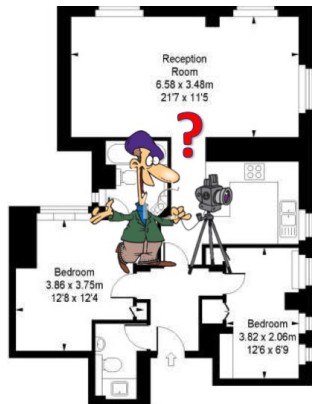
- ▷ Cabral & Furukawa, 2014
- ▷ Brualla et al., 2014
- Indoor localization (video, depth sensors)
 - Project Tango
 - SLAM work

Xiao & Furukawa, 2012

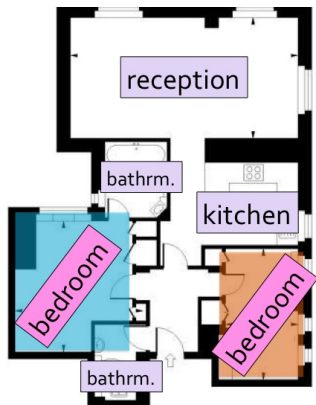


Cabral & Furukawa, 2014

Overview



Overview

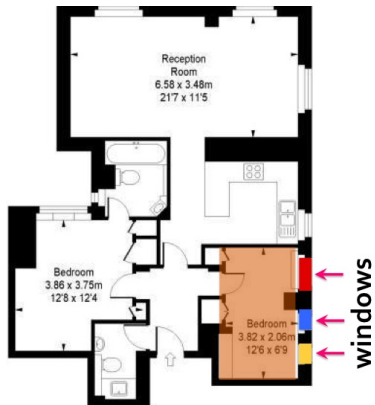


bedroom

Accurate **camera localization**:

- **Scene cues**

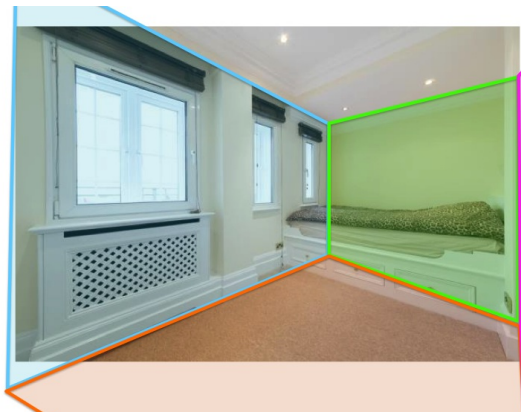
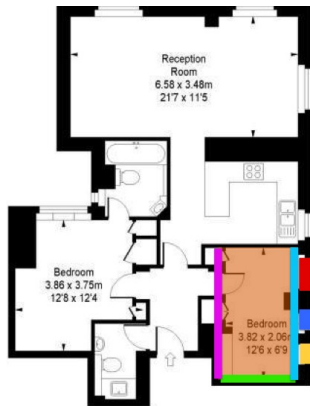
Overview



Accurate **camera localization**:

- **Scene cues**
- **Semantic cues**

Overview

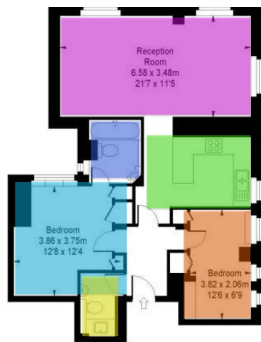


Accurate **camera localization**:

- **Scene cues**
- **Semantic cues**
- **Geometric cues** by exploiting the dimension information

Formulation

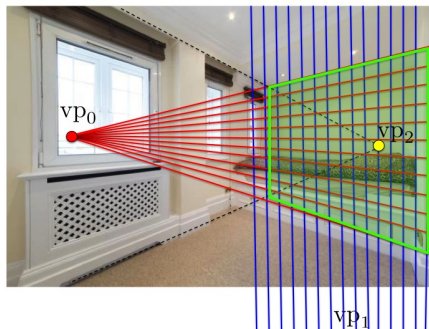
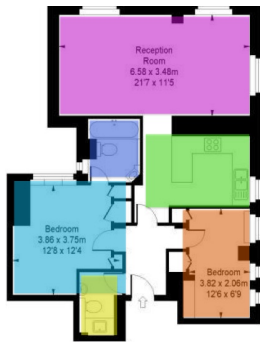
- $r \in \{1, \dots, R\}$... discrete random variable representing the room



Formulation

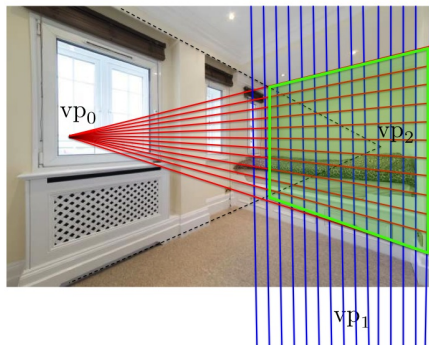
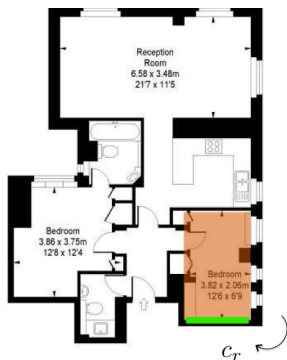
- $r \in \{1, \dots, R\}$... discrete random variable representing the room

Front wall is the plane defined by vp_0 and vp_1



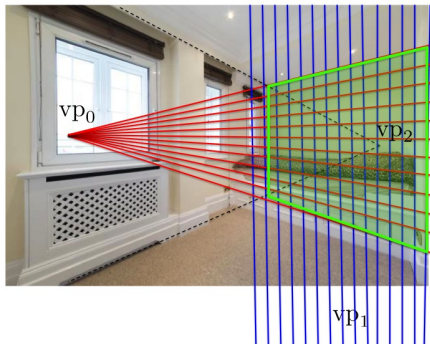
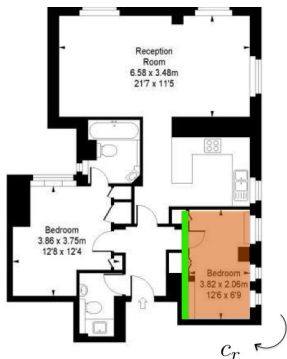
Formulation

- $r \in \{1, \dots, R\}$... discrete random variable representing the room
- $c_r \in \{1, \dots, |C_r|\}$... a discrete variable representing within room r which wall the picture is facing ($|C_r|$ the number of walls in a room)



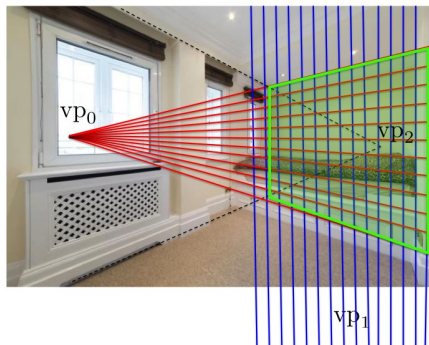
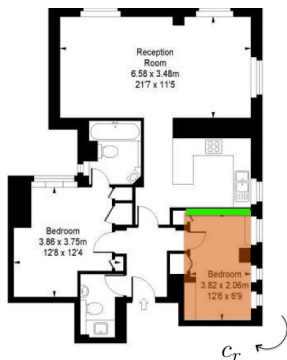
Formulation

- $r \in \{1, \dots, R\}$... discrete random variable representing the room
- $c_r \in \{1, \dots, |C_r|\}$... a discrete variable representing within room r which wall the picture is facing ($|C_r|$ the number of walls in a room)



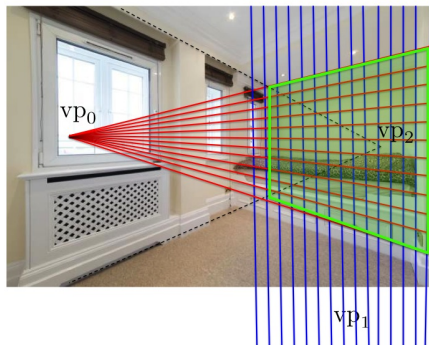
Formulation

- $r \in \{1, \dots, R\}$... discrete random variable representing the room
- $c_r \in \{1, \dots, |C_r|\}$... a discrete variable representing within room r which wall the picture is facing ($|C_r|$ the number of walls in a room)



Formulation

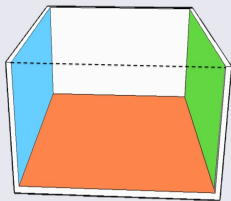
- $r \in \{1, \dots, R\}$... discrete random variable representing the room
- $c_r \in \{1, \dots, |C_r|\}$... a discrete variable representing within room r which wall the picture is facing ($|C_r|$ the number of walls in a room)



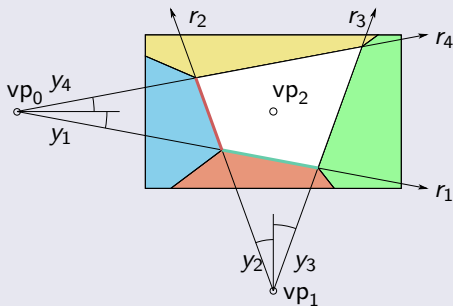
Formulation

- $r \in \{1, \dots, R\}$... discrete random variable representing the room
- $c_r \in \{1, \dots, |C_r|\}$... a discrete variable representing within room r which wall the picture is facing ($|C_r|$ the number of walls in a room)
- \mathbf{y} ... rays representing a room layout

Typical parametrization for room layout [Hedau et al., 2009]:



- Room is a 3D cuboid
- $\mathbf{y} = (y_1, y_2, y_3, y_4)$
- 4 rays needed to define it



Formulation

- $r \in \{1, \dots, R\}$... discrete random variable representing the room
- $c_r \in \{1, \dots, |C_r|\}$... a discrete variable representing within room r which wall the picture is facing ($|C_r|$ the number of walls in a room)
- \mathbf{y} ... rays representing a room layout
- We formulate the problem as inference in a Conditional Random Field with the following energy:

$$E(r, c_r, \mathbf{y}) = E_{scene_type}(r) + E_{layout}(r, c_r, \mathbf{y}) + E_{win}(r, c_r, \mathbf{y})$$

Energy Terms: Scene Type

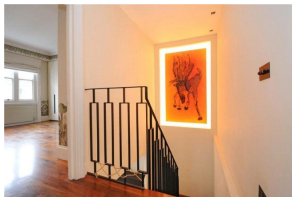
$$E(r, c_r, \mathbf{y}) = E_{\text{scene_type}}(r) + E_{\text{layout}}(r, c_r, \mathbf{y}) + E_{\text{win}}(r, c_r, \mathbf{y})$$

- **Potential:** Score of a scene classifier predicting scene type (e.g., bedroom, kitchen, reception)

Energy Terms: Scene Type

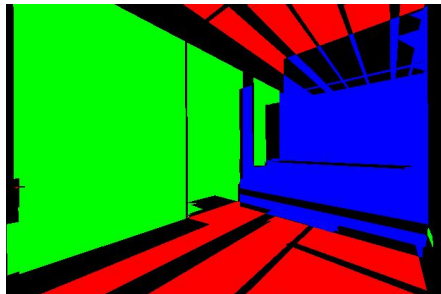
$$E(r, c_r, \mathbf{y}) = E_{\text{scene_type}}(r) + E_{\text{layout}}(r, c_r, \mathbf{y}) + E_{\text{win}}(r, c_r, \mathbf{y})$$

- **Potential:** Score of a scene classifier predicting scene type (e.g., bedroom, kitchen, reception)

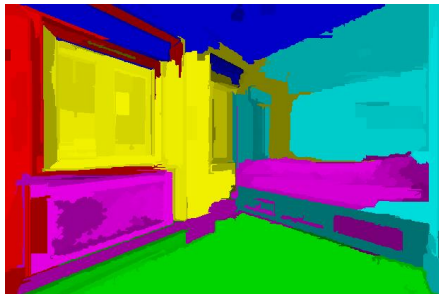


Energy Terms: Layout

$$E(r, c_r, \mathbf{y}) = E_{scene_type}(r) + E_{layout}(r, c_r, \mathbf{y}) + E_{win}(r, c_r, \mathbf{y})$$



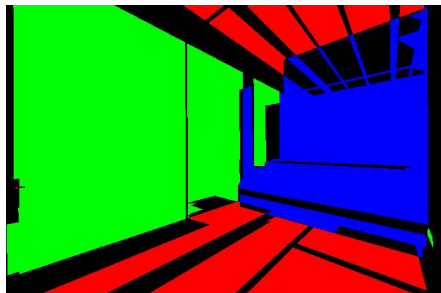
Orientation Map [Lee et al., 2009]



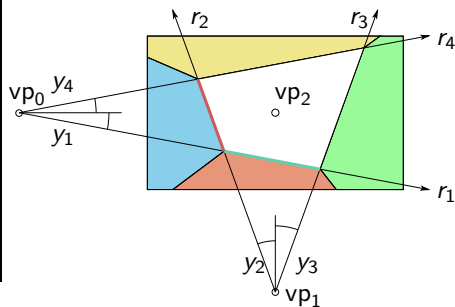
Geometric Context [Hedau et al., 2009]

Energy Terms: Layout

$$E(r, c_r, \mathbf{y}) = E_{scene_type}(r) + E_{layout}(r, c_r, \mathbf{y}) + E_{win}(r, c_r, \mathbf{y})$$



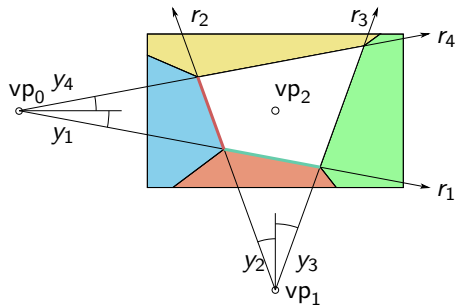
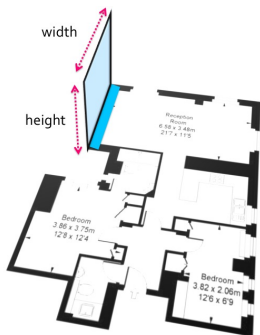
Orientation Map [Lee et al., 2009]



- **Potential:** Counts of blue, red, etc, pixels inside and outside of each wall
- Fast computation using *integral geometry* [Schwing et al., 2012]

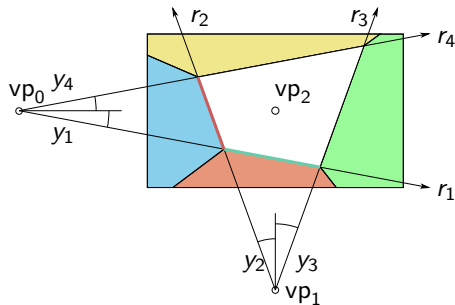
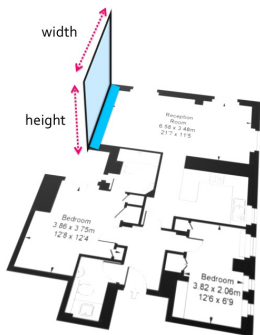
Energy Terms: Layout

$$E(r, c_r, \mathbf{y}) = E_{\text{scene_type}}(r) + E_{\text{layout}}(\boxed{r, c_r}, \mathbf{y}) + E_{\text{win}}(r, c_r, \mathbf{y})$$



Energy Terms: Layout

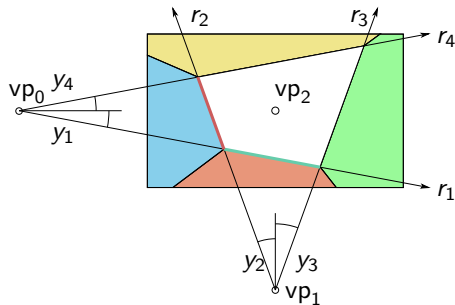
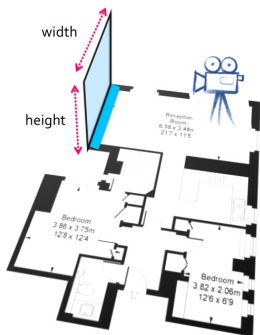
$$E(r, c_r, \mathbf{y}) = E_{\text{scene_type}}(r) + E_{\text{layout}}(\boxed{r, c_r}, \mathbf{y}) + E_{\text{win}}(r, c_r, \mathbf{y})$$



- $\mathbf{y} = (y_1, y_2, y_3, \cancel{y_4})$, $y_4 = f(r, c_r, y_1, y_2, y_3)$

Energy Terms: Layout

$$E(r, c_r, \mathbf{y}) = E_{scene_type}(r) + E_{layout}(r, c_r, \mathbf{y}) + E_{win}(r, c_r, \mathbf{y})$$



- $\mathbf{y} = (y_1, y_2, y_3, \cancel{y_4})$, $y_4 = f(r, c_r, y_1, y_2, y_3)$
- Additional constraint on \mathbf{y} : Camera is **inside** the room

Energy Terms: Windows

$$E(r, c_r, \mathbf{y}) = E_{scene_type}(r) + E_{layout}(r, c_r, \mathbf{y}) + E_{win}(r, c_r, \mathbf{y})$$

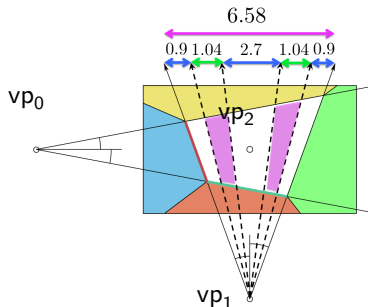
- Window-background segmentation



Energy Terms: Windows

$$E(r, c_r, \mathbf{y}) = E_{scene_type}(r) + E_{layout}(r, c_r, \mathbf{y}) + E_{win}(\boxed{r, c_r}, \mathbf{y})$$

- Window-background segmentation
- **Potential:** count window pixels inside and outside the window area



- We are minimizing the energy:

$$(r^*, c_r^*, \mathbf{y}^*) = \underset{r, c_r, \mathbf{y}}{\operatorname{argmin}} (E_{\text{scene_type}}(r) + E_{\text{layout}}(r, c_r, \mathbf{y}) + E_{\text{win}}(r, c_r, \mathbf{y}))$$

- We are minimizing the energy:

$$(r^*, c_r^*, \mathbf{y}^*) = \operatorname{argmin}_{r, c_r, \mathbf{y}} (E_{\text{scene_type}}(r) + E_{\text{layout}}(r, c_r, \mathbf{y}) + E_{\text{win}}(r, c_r, \mathbf{y}))$$

- Inference:
 - Exhaustive enumeration of r and c_r
 - Exact branch and bound inference for \mathbf{y} [Schwing & Urtasun, 2012]

- We are minimizing the energy:

$$(r^*, c_r^*, \mathbf{y}^*) = \operatorname{argmin}_{r, c_r, \mathbf{y}} (E_{\text{scene_type}}(r) + E_{\text{layout}}(r, c_r, \mathbf{y}) + E_{\text{win}}(r, c_r, \mathbf{y}))$$

- Inference:
 - Exhaustive enumeration of r and c_r
 - Exact branch and bound inference for \mathbf{y} [Schwing & Urtasun, 2012]
- We use S-SVM for training

Dataset

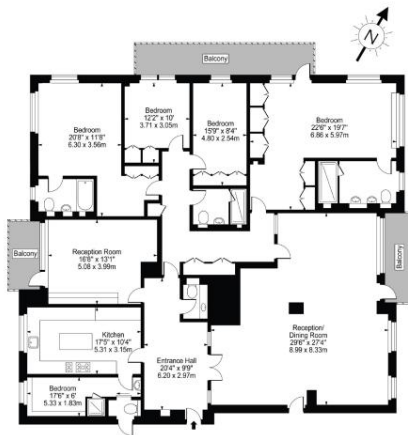
- We crawled a London apartment rental site

# apartments	215
# of images	1570
# of indoor images	1259
# images without GT alignment	82
avg. # rooms per apt	6
avg. # walls per apt	31
avg. # windows per apt	6
avg. # doors per apt	9



Apartments in Central London Are Not Small

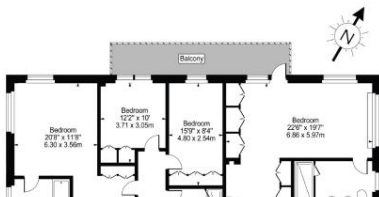
Approx. Gross Internal Area 2696 Sq Ft - 250.46 Sq M



Biggest apartment in dataset: 16 rooms, 5 bedrooms, 88 walls

Apartments in Central London Are Not Small

Approx. Gross Internal Area 2696 Sq Ft - 250.46 Sq M



Rent: £25,000 per month



Biggest apartment in dataset: 16 rooms, 5 bedrooms, 88 walls.

Results: Layout Estimation

- We assume we know which wall the camera is facing
- **Metrics:** Pixel accuracy for predicting 5 walls

	Layout error	Evaluations	Test time [s]
Schwing'12	13.88	16012.4	0.0208
Ours	11.81	1269.5	0.0019

Results: Layout Estimation

- We assume we know which wall the camera is facing
- **Metrics:** Pixel accuracy for predicting 5 walls

	Layout error	Evaluations	Test time [s]
Schwing'12	13.88	16012.4	0.0208
Ours	11.81	1269.5	0.0019

- 2% reduction in layout error

Results: Layout Estimation

- We assume we know which wall the camera is facing
- **Metrics:** Pixel accuracy for predicting 5 walls

	Layout error	Evaluations	Test time [s]
Schwing'12	13.88	16012.4	0.0208
Ours	11.81	1269.5	0.0019

- 2% reduction in layout error
- 10 times less branching operations

Results: Layout Estimation

- We assume we know which wall the camera is facing
- **Metrics:** Pixel accuracy for predicting 5 walls

	Layout error	Evaluations	Test time [s]
Schwing'12	13.88	16012.4	0.0208
Ours	11.81	1269.5	0.0019

- 2% reduction in layout error
- 10 times less branching operations
- 10x speedup

Results: Camera Localization

- **Metrics:** % of correct assignments of front wall to the apartment wall

	Aspect	+Scene	+Room
Random	0.0328	0.1138	0.1954
Ours (no windows)	0.0686	0.1945	0.2654
Ours (windowGT)	0.2128	0.4737	0.5995
Ours (window)	0.1670	0.3982	0.5080

Results: Camera Localization

- **Metrics:** % of correct assignments of front wall to the apartment wall

	Aspect	+Scene	+Room
Random	0.0328	0.1138	0.1954
Ours (no windows)	0.0686	0.1945	0.2654
Ours (windowGT)	0.2128	0.4737	0.5995
Ours (window)	0.1670	0.3982	0.5080

Aspect: Only aspect ratio information (and not scene) used

Results: Camera Localization

- **Metrics:** % of correct assignments of front wall to the apartment wall

	Aspect	+Scene	+Room
Random	0.0328	0.1138	0.1954
Ours (no windows)	0.0686	0.1945	0.2654
Ours (windowGT)	0.2128	0.4737	0.5995
Ours (window)	0.1670	0.3982	0.5080

+*Scene*: Aspect information and scene classifier are used

Results: Camera Localization

- **Metrics:** % of correct assignments of front wall to the apartment wall

	Aspect	+Scene	+Room
Random	0.0328	0.1138	0.1954
Ours (no windows)	0.0686	0.1945	0.2654
Ours (windowGT)	0.2128	0.4737	0.5995
Ours (window)	0.1670	0.3982	0.5080

+*Room*: We know which room the picture was taken in

Results: Camera Localization

- **Metrics:** % of correct assignments of front wall to the apartment wall

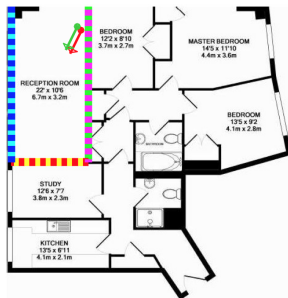
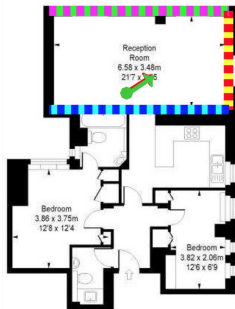
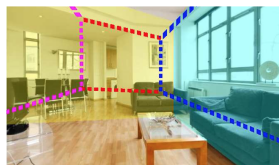
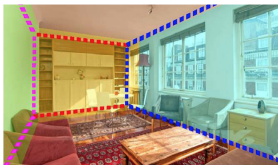
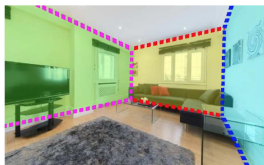
	Aspect	+Scene	+Room
Random	0.0328	0.1138	0.1954
Ours (no windows)	0.0686	0.1945	0.2654
Ours (windowGT)	0.2128	0.4737	0.5995
Ours (window)	0.1670	0.3982	0.5080

Results: Camera Localization

- **Metrics:** % of correct assignments of front wall to the apartment wall

	Aspect	+Scene	+Room
Random	0.0328	0.1138	0.1954
Ours (no windows)	0.0686	0.1945	0.2654
Ours (windowGT)	0.2128	0.4737	0.5995
Ours (window)	0.1670	0.3982	0.5080

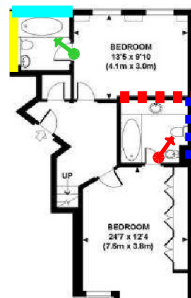
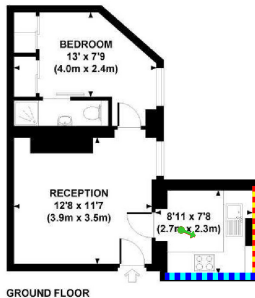
Results: Joint Layout and Localization



Red arrow: Groundtruth camera

Green arrow: Predicted camera

Results: Joint Layout and Localization

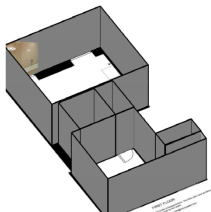


Red arrow: Groundtruth camera

Green arrow: Predicted camera

Results: Reconstruction

Window+Aspect



1 images out of 4
2 walls out of 8

+Scene



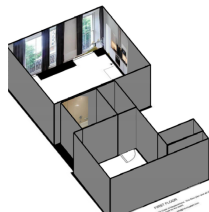
4 images out of 4
8 walls out of 8

+Room



4 images out of 4
8 walls out of 8

Ground-truth



-
-



HDVA



HDVA

Summary

- Problem of apartment 3D reconstruction from monocular imagery
- Model that jointly solves for localization and room layout estimation by exploiting floor-plans
- Real-time inference
- Results:
 - We improve layout prediction over past work
 - Achieve good localization performance
- Dataset with 215 apartments and all annotations available:

<http://www.cs.toronto.edu/~fidler/projects/rent3D.html>

Alex on the Market Next Year



Next year

Thank You
Welcome to our poster at #9!