

Super Resolution for Automated Target Recognition

Maksim Levental

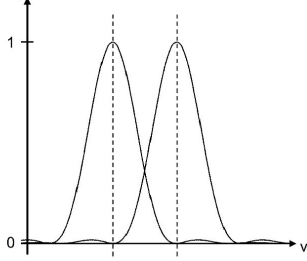


Fig. 1. Rayleigh's criterion[2]

Abstract—Super resolution is the process of producing high-resolution images from low-resolution images while preserving ground truth about the subject matter of the images and potentially inferring more such truth. Algorithms that successfully carry out such a process are broadly useful in all circumstances where HR imagery is either difficult or impossible to obtain. In particular we look towards super resolving images collected using longwave infrared cameras since high resolution sensors for such cameras do not currently exist. We present an exposition of motivations and concepts of super resolution in general and current techniques, with a qualitative comparison of such techniques. Finally we suggest directions for future research.

1 INTRODUCTION

Super-resolution (SR) is a collection of methods¹ that augment the resolving power of an imaging system. Here, and in the forthcoming, by resolving power we mean the ability of an imaging device to distinguish distinct but proximal objects in the scene. If such objects are modeled as point sources of light then the resolving power of the imaging system is defined by Rayleigh's criterion: two point sources are considered *resolved* when the first diffraction maximum² of one point source (at most) coincides with the first minimum of the other (see figure 1).

SR techniques yield high-resolution (HR) images from one or more observed low-resolution (LR) images by restoring lost fine details and reversing degradations produced by imperfect imaging systems. In the case when a single LR source image is used to construct the HR correspondent, the techniques

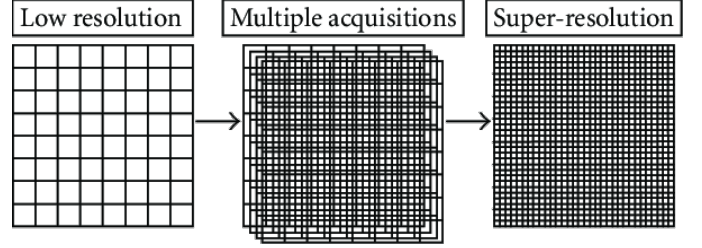


Fig. 2. Multiple image super resolution[3]

are referred to as single-image-super-resolution (SISR) techniques. These techniques typically operate by either learning some mapping from low resolution chips (uniform partitions of the image, e.g. 3×3 pixels) to higher resolution chips that are highly similar (according to some metrics) and which obey regularity constraints (e.g. agreement at edges). In the case when multiple LR source images are used to construct the single HR correspondent, the techniques are referred to as multiple-image-super-resolution (MISR) techniques. MISR techniques rely on non-redundant and yet pertinent information in multiple images of the same scene (see figure 2). Note that for such information to exist there should be sub-pixel³ shifts in either the imaging system or the scene between consecutive images.

For typical imaging use-cases, high resolution images are preferable to low resolution images; higher resolutions are desirable in and of themselves and as inputs to later image processing transformations that can degrade image quality (e.g. by virtue of quantization or compression). In theory the resolving power of an imaging system is primarily determined by the number of independent sensor elements that comprise that imaging system (each of which collects a component of the ultimate image). Naturally then, a way to increase the resolution of such a system is to increase the density of such sensor elements per unit area. Unfortunately, and counter-intuitively, since the number of photons incident on each sensor decreases as the sensor shrinks, shot noise⁴ thwarts that idea. Furthermore, while sensor density is primary, secondary effects due to optics limit resolution as well; the point spread of a lens (distortion of a point source due to diffraction), chromatic aberrations (distortion due to differing indices of refraction for differing wavelengths of light), and motion blur all function to obscure or erase details from the image.

In domains such as satellite/aerial photography, medical imaging, and facial recognition, high-resolution reconstruction of low-resolution samples is eminently useful since ab-initio

¹We will often use the verb form "to super resolve" in order to denote the use of one or more such methods.

²The amplitude of the diffraction pattern (known as the Airy pattern) of a monochromatic point source through a circular aperture is given by

$$I(\theta) = I_0 \left[\frac{2J_1(ka \sin \theta)}{ka \sin \theta} \right]^2$$

where I_0 is peak intensity (at the center), $k = \frac{2\pi}{\lambda}$ is the wave number of the light, θ is the angle of observation, and J_1 is the Bessel function of the first kind of order one[1]. It is maxima/minima of this function that Rayleigh's criterion concerns.

³For example when a point source wholly captured by one sensor element shifts to distributing energy equally amongst the same element and a direct adjacent.

⁴TODO

acquisition of high-resolution images is either logistically difficult or impossible due to aforementioned imaging apparatus limitations. For example in the instance of satellite imagery, acquisition of high-resolution imagery is primarily hampered by optics and physics⁵. In contrast, in the cases of medical imaging (where procedures are invasive and patient exposure time needs to be minimized[6]) and facial recognition (e.g. for purposes of surveillance) the primary challenge is logistics and access to repeat collection opportunities.

The benefits of enhancing images using SR techniques include not only more pleasing or more readily interpretable images for human consumption but higher quality inputs for automated learning systems as well. In particular object detection systems trained on super-resolved images outperform those trained on the low resolution originals[7]. Indeed this is our ultimate goal - not super-resolution per se but super-resolution in the service of improved object detection performance for longwave-infrared (LWIR) imagery. Note that while practically speaking, there exist hardware and software solutions for increasing the resolution of an imaging system, we, owing to a "Ship of Theseus" consideration, discount such propositions. We instead take low resolution images as given and seek techniques that allow for ex post facto reconstruction or inference of precise details. This necessarily constrains techniques under consideration to be algorithmic in nature and software in practice.

The rest of this survey is outlined as follows: Section 2 introduces notation and the model of imaging that will be the mathematical framework for the proceeding sections, Section 3 surveys classical techniques (those that do not employ neural networks), Section 4 surveys neural-network techniques with heavy emphasis on deep learning (i.e. deep networks), Section 5 discusses the scope and goals of the author's research program, and Section 6 summarizes.

2 BACKGROUND

2.1 Imaging systems

An imaging sensor is a device that converts an optical image into a digital signal. Charge-coupled devices (CCD) and complementary metal-oxide-semiconductor (CMOS) devices are the most common imaging sensors; CCDs have better performance while CMOS devices are newer and less expensive. A third type that's of particular interest to us is the microbolometer, which is used as a detector in thermal cameras.

CCDs consist of densely packed two-dimensional arrays of buried channel⁶ MOS capacitors (see figure 3) with an

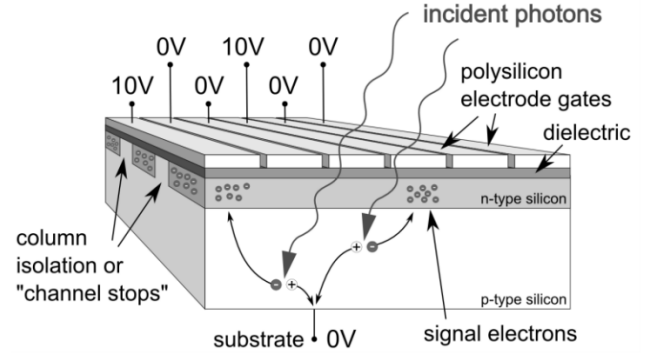


Fig. 3. CCD buried channel MOS capacitor[8]

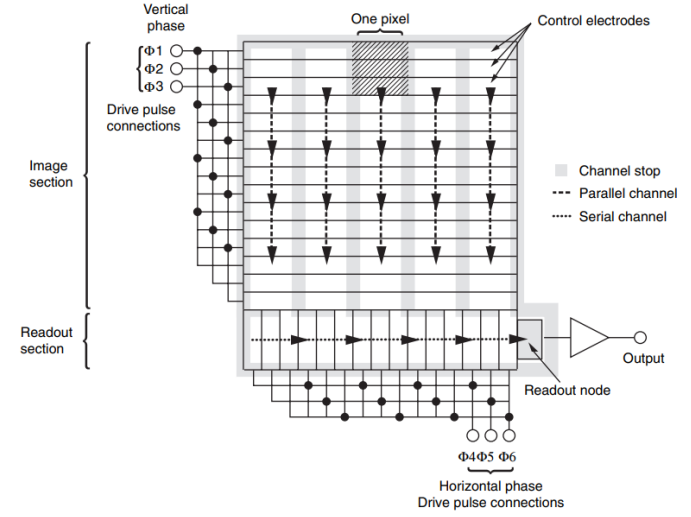


Fig. 4. CCD array[11]

individual MOS capacitor being the fundamental photon detecting element. Individual MOS capacitors are biased by a gate voltage such that a signal electron trap (potential well) is produced in the n-type silicon (referred to as the n-channel). When the array is illuminated some of the light is reflected, some is absorbed by the gate, and the rest if transmitted into the silicon. While the absorption coefficients of poly-silicon are fairly high ($4 \times 10^4 \text{cm}^{-1}$ at 450nm and $1.2 \times 10^4 \text{cm}^{-1}$ at 550nm) alternative electrode materials exist with better transmission coefficients[9]. The photons that are transmitted interact with electrons in the valence band of the silicon and excite them into the conduction band, thereby creating electron-hole pairs that either scatter or recombine. For high-quality silicon, the lifetime of such a pair is several milliseconds (before recombination)[10]. Of those that scatter the electrons move by diffusion into the electron trap while the holes migrate to the substrate (i.e. out of the capacitor).

CCD arrays consist of two sub-arrays: and image section and a readout section. The image section is arranged with every third stripe of electrode tied electrically to form three sets of equipotentials. In figure 4 these equipotentials are labeled $\Phi1, \Phi2, \Phi3$, and taken together constitute a vertical register (VR). They function to move the collected charge down one line at a time, using charge coupling, while the

⁵Rayleigh's criterion implies that the angular resolution R of a telescope with optical diameter $D = 2.4\text{m}$ observing visible light ($\sim 500\text{nm}$) is approximately[4]

$$R \approx 1.220 \frac{\lambda}{D} = 1.220 \frac{500\text{nm}}{2.4\text{m}} \approx 0.06\text{arcsec}$$

From an altitude of 250 km this corresponds to a ground sample distance of 6cm. This loss of resolving power is further exacerbated by refraction through turbulent atmosphere[5].

⁶Buried channel as distinct from surface channel. In surface channel MOS capacitors signal charge is stored at the Si-SiO₂ interface, which can lead to charge trapping during the charge transfer process.

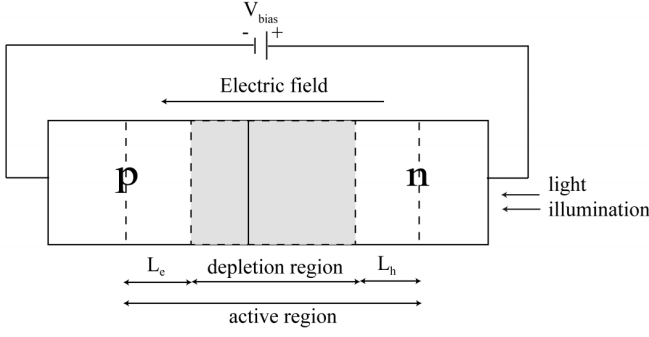


Fig. 5. Photodiode schematic. L_e , L_h are electron, hole diffusion lengths respectively[12]

channel stops function to prevent diffusion of charge across channels. The VR mechanism that moves collected charge functions like this:

- 1) Suppose initially there's a collection of charge on each channel at Φ_1 and only Φ_1 . Note this means Φ_2 is at 0v (again just as in figure 3).
- 2) Φ_2 is positively biased to 10V. This diffuses the collection of charge under both Φ_1 and Φ_2 .
- 3) Φ_1 is set to 0v. This concentrates the collection of charge under Φ_2 .
- 4) The same process is repeated with Φ_2 , Φ_3 and Φ_3 , Φ_1 .
- 5) Goto step 1.

At the bottom of the image section Φ_3 transfers all signal charges to the horizontal register (HR) which functions much like the VR except much faster: the HR must transfer every line independently of all other lines to the read-out node. For a given image the first pixel read at the read-out node is the bottom-right pixel and the last pixel read is the top-left pixel.

CMOS sensors consist of arrays of photodiodes (see figure 5). A photodiode is a p-n junction⁷ operated in reverse bias mode⁸. When a photon of sufficient energy is absorbed by the diode, it creates an electron-hole pair. If the creation event happens within the active region then the hole moves out through the p-type material and the electron moves out through the n-type material. This establishes a *photocurrent* that can be read by a reading circuit (see figure 6). CMOS sensor arrays do not shift the charge from row to row. In a CMOS sensor array, each pixel contains a transistor M_{sel} controlled by the voltages applied across a row (see figure 7). In order to read one row of pixels, a row line is raised high to turn on (close) all the M_{sel} transistors in the row. This brings the signals from all the pixels in that row to the shifter register below by way of the column lines. The shift register then outputs the values of the pixels.

A microbolometer is a specific type of bolometer used as a detector in a thermal camera. Infrared radiation with wavelengths between $7.5\text{--}14\text{ }\mu\text{m}$ strikes the detector material,

⁷The interface between a p-type semiconductor (excess holes, i.e. positive charge carriers) and an n-type (excess electrons, i.e. negative charge carriers) semiconductor.

⁸With the p-type material at a lower voltage than the n-type. This causes both the holes and the electrons to flow away from the junction creating a depletion zone.

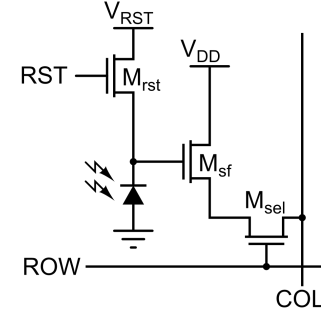


Fig. 6. Three transistor "pixel". M_{rst} is the reset transistor (enabling the photodiode to dump charge), M_{sf} buffers the charge on the photodiode (so that it can be read without loss), and M_{sel} enables a whole row of pixels to be read simultaneously (since all pixels in a physical row are tied to the same row line).

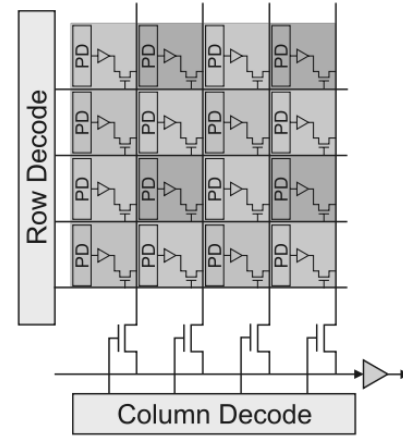


Fig. 7. CMOS array

heating it, and thus changing its electrical resistance (see figure 8). This resistance change is measured and processed into temperatures which can be used to create an image.

2.2 Mathematical notation

Upper case plain latin X, Y denote channel \times row \times column *tensors*⁹ representing LR and HR images respectively,

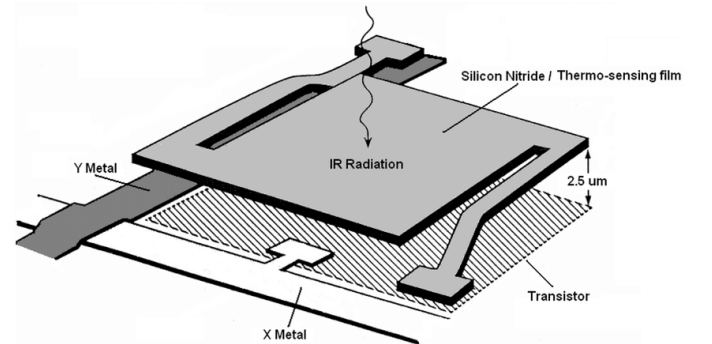


Fig. 8. Microbolometer[13]

⁹A multidimensional array[14]. Not to be confused with the algebraic object.

with $(0, 0)$ corresponding to the top left corner of the image. Often for the sake of simplicity we consider greyscale images in which case we omit the channel dimension. Lower case plain latin x, y denote LR and HR patches¹⁰ respectively. D, H, F, G variously refer to functions that operate on images. Bolded latin X, Y denotes batches.

2.3 Imaging model

Figure 9 shows a conceptual model of the imaging process as carried out by an imaging system. The input to the system is a natural scene that is in effect sampled by the imaging system. In the idealized case the sampling is done at (or above) the Nyquist rate and no aliasing occurs. In practice there is noise and loss introduced at every step of the process: atmospheric turbulence plays a role at large distances, motion produces multiple views of the same scene but also induces blur, imperfections of the lenses further blur the image, and finally down-sampling by the sensor elements into pixels produces aliasing artifacts¹¹. The noisy, blurry, down-sampled images are then further degraded by sensor noise. Each such image we call an LR sample.

Let Y denote an idealized HR image of the scene from some fixed vantage point and assume the imaging system collects K LR samples X_k of Y . Formally the X_k are related to Y by

$$X_k = (D_k \circ H_k \circ A_k)(Y) + \epsilon \quad (1)$$

where for the k th sample A_k is the affine transformation representing motion (rigid and perspective shift), H_k represents the composite blur operator (motion and optics blur), D_k represents the down-sampling operator, and ϵ represents the composite noise (environment and sensor noise). In general A_k, H_k, D_k are highly degenerate functions for which the corresponding inverse problems are ill-posed without regularization and conditioning.

3 CLASSICAL ALGORITHMS

4 DEEP LEARNING ALGORITHMS

5 FUTURE RESEARCH

6 CONCLUSION

7 APPENDIX

TODO: work out diffraction circular aperture TODO: work-out poisson noise

ACKNOWLEDGMENTS

REFERENCES

- [1] J. W. Goodman, *Introduction to Fourier Optics*, 3rd. Roberts & Co. Publishers, 2005, pp. 76–78.
- [2] P. Scholz, “Focused ion beam created refractive and diffractive lens techniques for the improvement of optical imaging through silicon,” PhD thesis, Jul. 2012. DOI: 10.14279/depositonce-3270.
- [3] J. Kennedy, O. Israel, A. Frenkel, R. bar-shalom, and H. Azhari, “Improved image fusion in pet/ct using hybrid image reconstruction and super-resolution,” *International journal of biomedical imaging*, vol. 2007, p. 46846, Jan. 2007. DOI: 10.1155/2007/46846.
- [4] L. R. F.R.S., “Xxxi. investigations in optics, with special reference to the spectroscope,” *The London, Edinburgh, and Dublin Philosophical Magazine and Journal of Science*, vol. 8, no. 49, pp. 261–274, 1879. DOI: 10.1080/14786447908639684. eprint: <https://doi.org/10.1080/14786447908639684>. [Online]. Available: <https://doi.org/10.1080/14786447908639684>.
- [5] D. L. Fried, “Optical resolution through a randomly inhomogeneous medium for very long and very short exposures,” *J. Opt. Soc. Am.*, vol. 56, no. 10, pp. 1372–1379, Oct. 1966. DOI: 10.1364/JOSA.56.001372. [Online]. Available: <http://www.osapublishing.org/abstract.cfm?URI=josa-56-10-1372>.
- [6] E. Van Reeth, I. W. K. Tham, C. H. Tan, and C. L. Poh, “Super-resolution in magnetic resonance imaging: A review,” *Concepts in Magnetic Resonance Part A*, vol. 40A, no. 6, pp. 306–325, 2012. DOI: 10.1002/cmr.a.21249. eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1002/cmr.a.21249>. [Online]. Available: <https://onlinelibrary.wiley.com/doi/abs/10.1002/cmr.a.21249>.
- [7] J. Shermeyer and A. V. Etten, “The effects of super-resolution on object detection performance in satellite imagery,” *CoRR*, vol. abs/1812.04098, 2018. arXiv: 1812.04098. [Online]. Available: <http://arxiv.org/abs/1812.04098>.
- [8] M. Robbins, “Final test guideline,” May 2014.
- [9] M. Bass, C. DeCusatis, J. Enoch, V. Lakshminarayanan, G. Li, C. Macdonald, V. Mahajan, and E. Van Stryland, *Handbook of Optics, Third Edition Volume I: Geometrical and Physical Optics, Polarized Light, Components and Instruments(Set)*, 3rd ed. New York, NY, USA: McGraw-Hill, Inc., 2010, ISBN: 0071498893, 9780071498890.
- [10] J. R. Janesick, T. Elliott, S. Collins, M. M. Blouke, and J. Freeman, “Scientific Charge-Coupled Devices,” *Optical Engineering*, vol. 26, no. 8, pp. 692–714, 1987. DOI: 10.1117/12.7974139. [Online]. Available: <https://doi.org/10.1117/12.7974139>.
- [11] J. Pawley, *Handbook of Biological Confocal Microscopy*, ser. Cognition and Language. Springer, 1995, pp. 918–919, ISBN: 9780306448263. [Online]. Available: <https://books.google.com/books?id=16Ft5k8RC-AC>.
- [12] Y. Xu, “Fundamental characteristics of a pinned photodiode cmos pixels,” 2015.

¹⁰ $k \times k$ pixel window, e.g. 3×3 .

¹¹TODO

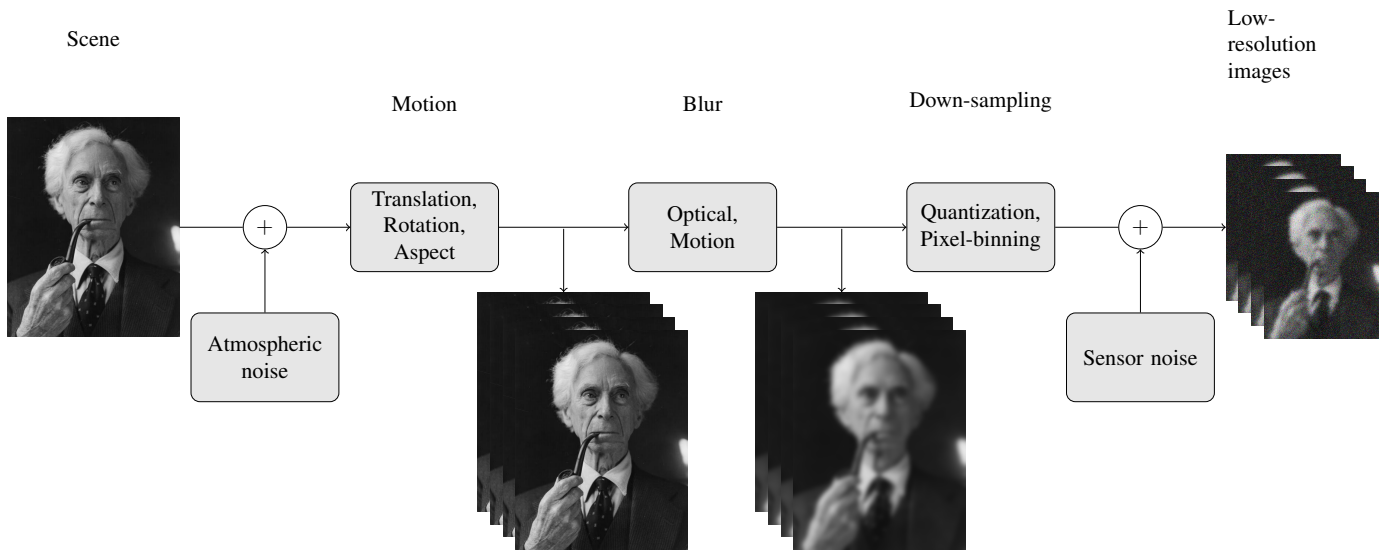


Fig. 9. The imaging model illustrating the relationship between a scene and final low-resolution images due to noise, motion, blur, and sampling.

- [13] R. Ambrosio, M. Moreno, J. Mireles Jr., A. Torres, A. Kosarev, and A. Heredia, "An overview of uncooled infrared sensors technology based on amorphous silicon and silicon germanium alloys," *physica status solidi c*, vol. 7, no. 3-4, pp. 1180–1183, 2010. DOI: 10.1002/pssc.200982781. eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1002/pssc.200982781>. [Online]. Available: <https://onlinelibrary.wiley.com/doi/abs/10.1002/pssc.200982781>.
- [14] I. Goodfellow, Y. Bengio, and A. Courville, *Deep Learning*. The MIT Press, 2016, p. 33, ISBN: 0262035618, 9780262035613.