

Supplementary Information:
**Plasmonic antenna coupling to hyperbolic phonon-polaritons for
sensitive and fast mid-infrared photodetection with graphene**

S. Castilla et al.

SUPPLEMENTARY NOTE 1: OPTICAL MODELING

We perform optical calculations with the full vector 3D finite-difference time domain (FDTD) method using Lumerical software. The computational cell dimensions are $6 \times 6 \times 6 \text{ } \mu\text{m}^2$ with perfectly matched layer (PML) conditions employed on all boundaries. We use a grid of 5 nm in lateral (x,y) and 2 nm in vertical (z) directions. The light source is modeled as an incident plane wave reaching the device through an open aperture of $5.9 \times 5.9 \text{ } \mu\text{m}^2$ with spectral range 5.5-10.5 μm . The dimensions of the device layers are described in the *Fabrication* section in Methods.

We fit the Au dielectric function using the Drude model $\varepsilon(\omega) = \varepsilon_\infty - \omega_p^2/(\omega^2 + i\Gamma\omega)$ with ω_p the plasma frequency and Γ the plasma collision rate. hBN is optically anisotropic¹, with different permittivities along the in-plane ($\perp c$) and out-of-plane ($\parallel c$) directions. We fit both using the Lorentz model $\varepsilon(\omega) = \varepsilon_\infty + s\omega_0^2/(\omega_0^2 - \omega^2 - i\gamma\omega)$, where s is a dimensionless coupling factor, ω_0 the normal frequency of vibration and γ the decay rate amplitude. Supplementary Table 1 shows the parameters for Au and hBN. The refractive indices used for SiO_2 and Al_2O_3 are taken from literature², while for Si we use $n_{\text{Si}} = 3.42$. Graphene is implemented as a 2D surface with optical response modeled by the Kubo conductance³ $\sigma = \sigma_{\text{intra}} + \sigma_{\text{inter}}$, where:

$$\sigma_{\text{intra}} = \frac{ie^2}{\pi\hbar^2\Omega} \int_0^\infty \epsilon(\partial_\epsilon f(-\epsilon; \mu, T_e) - \partial_\epsilon f(\epsilon; \mu, T_e)) d\epsilon \quad (1)$$

$$\sigma_{\text{inter}} = \frac{ie^2\Omega}{\pi\hbar^2} \int_0^\infty \left(\frac{f(-\epsilon; \mu, T_e) - f(\epsilon; \mu, T_e)}{\Omega^2 - 4(\epsilon/\hbar)^2} \right) d\epsilon \quad (2)$$

Here Ω is defined as $\Omega = \omega + i\tau_{\text{opt}}^{-1}$, $\tau_{\text{opt}} = 200 \text{ fs}$ is the assumed electron relaxation time, $f(\epsilon; \mu, T_e) = [e^{(\epsilon-\mu)/k_B T_e} + 1]^{-1}$ is the Fermi-Dirac distribution and $\partial_\epsilon = \partial/\partial\epsilon$.

Material	ε_∞	s	ω_0, ω_p [eV]	γ_0, Γ [eV]
hBN ¹ $\perp c$	4.87	1.83	0.17	0.87
hBN ¹ $\parallel c$	2.95	0.61	0.0925	0.25
Au ⁴	10.78	-	9.13	0.07

Supplementary Table I: Dielectric permittivity parameters of Au and hBN.

SUPPLEMENTARY NOTE 2: THERMOELECTRIC MODELING

We assume the quasi-continuous-wave case and solve the heat dissipation equation⁵:

$$-\nabla \cdot (\kappa \nabla T_e) = \nabla \Pi \cdot \mathbf{j}_q - \tau_{\text{e-ph}}^{-1} c_e \Delta T + \alpha P_{\text{dens}} \quad (3)$$

where κ is the electronic thermal conductivity, T_e the electronic temperature, $\Pi = ST_e$ the Peltier coefficient, S the Seebeck coefficient, $\mathbf{j}_q = -\sigma S \nabla T_e$ the local thermoelectric current, σ the electrical conductivity, $\tau_{\text{e-ph}}$ the average cooling time (3 ps), c_e the electronic heat capacity, $\Delta T = T_e - T_l$, T_l is the lattice temperature, α the absorption fraction and the incident power density $P_{\text{dens}} = P_{\text{in}}/A_{\text{diff}}$, where P_{in} is the source power and $A_{\text{diff}} = \lambda^2/\pi$ is the diffraction-limited area. The term $c_e/\tau_{\text{e-ph}}$ in Supplementary Equation 3 is equivalent to Γ_{cool} , where Γ_{cool} ($\sim 4\text{-}5 \times 10^4 \text{ W/Km}^2$) is the interfacial heat conductivity^{6,7}. Due to the large lattice heat capacity (compared to the electronic one) we assume constant $T_l = 300 \text{ K}$. The graphene parameters $\sigma, S, c_e, \kappa, \Gamma_{\text{cool}}$ and α are functions of position, depending on both the local Fermi level E_F and the local temperature T_e , making this a highly non-linear problem. The self-consistent solution of Supplementary Equation 3 provides the T_e distribution from which the thermoelectric voltage is obtained:

$$V_{\text{PTE}} = W^{-1} \int_0^W \int_0^L S \nabla T_e dx dy \quad (4)$$

where the length L is the distance between the contacts (assumed in the x direction, as shown in Fig. 1a-b in the main text) and W the width of the graphene channel (y direction, as shown in Fig. 1a-b). The photocurrent

is then $I_{\text{PTE}} = V_{\text{PTE}}/R_{\text{D}}$, where $R_{\text{D}} = R_{\text{g}} + 2R_{\text{c}}$ with R_{g} is the resistance of graphene channel $R_{\text{g}} = \int_0^L \sigma^{-1}(x)dx$, $\sigma(x) = \int_0^W \sigma(x, y)dy$ and R_{c} the contact resistance of the device. The responsivity and NEP are calculated following the description in *Responsivity and NEP calculation* section in Methods.

Most parameters in the above equations depend on the local Fermi level and electronic temperature. The former is obtained from the graphene charge density $E_{\text{F}} = \hbar v_{\text{F}} \sqrt{\pi n^{(0)}}$, where $n^{(0)}$ is calculated by electrostatic simulations (see Supplementary Note 3 and Supplementary Figure 4) using the ratio $\epsilon_{\text{hBN}}/d_{\text{hBN}}$ as determined by fitting the measured device resistance (see Supplementary Note 4). At finite temperature we obtain the chemical potential from the solution of $\int_0^{\infty} v(\epsilon) f(\epsilon; \mu, T_e) - f(\epsilon; -\mu, T_e) d\epsilon = E_{\text{F}}^2 / \pi \hbar^2 v_{\text{F}}^2$, where $v(\epsilon) = \frac{2|\epsilon|}{\pi \hbar^2 v_{\text{F}}^2}$ is the graphene density of states at energy ϵ and $v_{\text{F}} = 1 \times 10^6$ m/s the graphene Fermi velocity. The rest of the graphene electrical and thermal parameters are calculated as follows:

Electrical conductivity:

$\sigma(\mu, T_e) = \int_{-\infty}^{\infty} \sigma(\epsilon) \partial_{\epsilon} f(\epsilon; \mu, T_e) d\epsilon$, where $\sigma(\epsilon) = q[\mu_{\text{q}} n(\epsilon) + \bar{\mu} n^*(\epsilon)]$, with charge carrier mobility $\mu_{\text{q}} = \mu_{\text{e}}(\mu_{\text{h}})$ for $\epsilon > 0$ ($\epsilon < 0$) and $\bar{\mu} = (\mu_{\text{e}} + \mu_{\text{h}})/2$. The effective residual local charge fluctuation at energy ϵ is assumed to be $n^*(\epsilon) = \sqrt{n(\epsilon)^2 + n_0^{*2}} - n(\epsilon)$, where $n(\epsilon) = \frac{e^2}{\pi \hbar^2 v_{\text{F}}^2}$ is the graphene charge density at energy ϵ and n_0^* is the residual local charge fluctuations in the charge neutrality point^{8,9}.

Seebeck coefficient is given by the general Mott formula¹⁰:

$$S(\mu, T_e) = -(|e|T_e \sigma)^{-1} \int_{-\infty}^{\infty} (\epsilon - \mu) \sigma(\epsilon) \partial_{\epsilon} f(\epsilon; \mu, T_e) d\epsilon$$

Thermal capacity¹¹:

$$c_{\text{e}}(\mu, T_e) = \partial_{T_e} \int_0^{\infty} (v(\epsilon) \epsilon [f(\epsilon; \mu, T_e) + f(\epsilon; -\mu, T_e)]) d\epsilon$$

Thermal conductivity: is given by the Wiedemann-Franz law $\kappa(\mu, T_e) = L_0 \sigma(\mu, T_e) T_e$, where $L_0 = 2.44 \times 10^{-8} \text{W}\Omega\text{K}^{-2}$ is the Lorenz number.

SUPPLEMENTARY NOTE 3: ELECTROSTATIC MODELING

The surface charge density of the graphene sheet is calculated by solving the Poisson equation after applying the appropriate voltages V_{L} and V_{R} at the two branches of the split gate. The graphene channel is introduced as a grounded surface above the hBN dielectric spacer layer. The nearby metal contacts are also set to ground. The dielectric constants of Si, SiO₂ and hBN are set to 11.7, 3.9 and 3.5¹²⁻¹⁴ respectively. The surface charge density is calculated for both symmetric and anti-symmetric gating cases (see Supplementary Figure 4a).

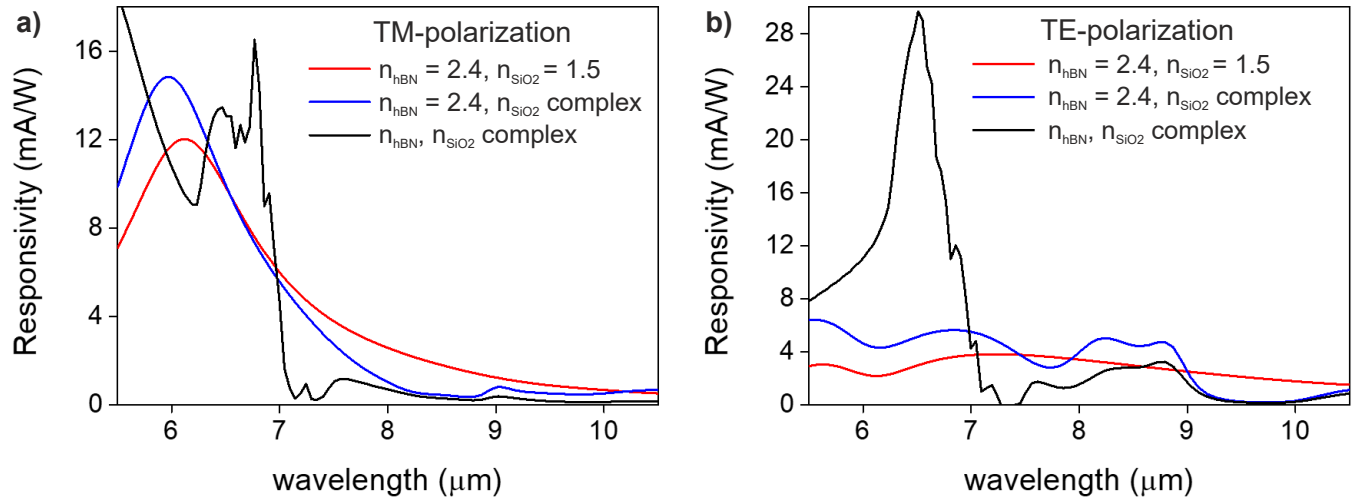
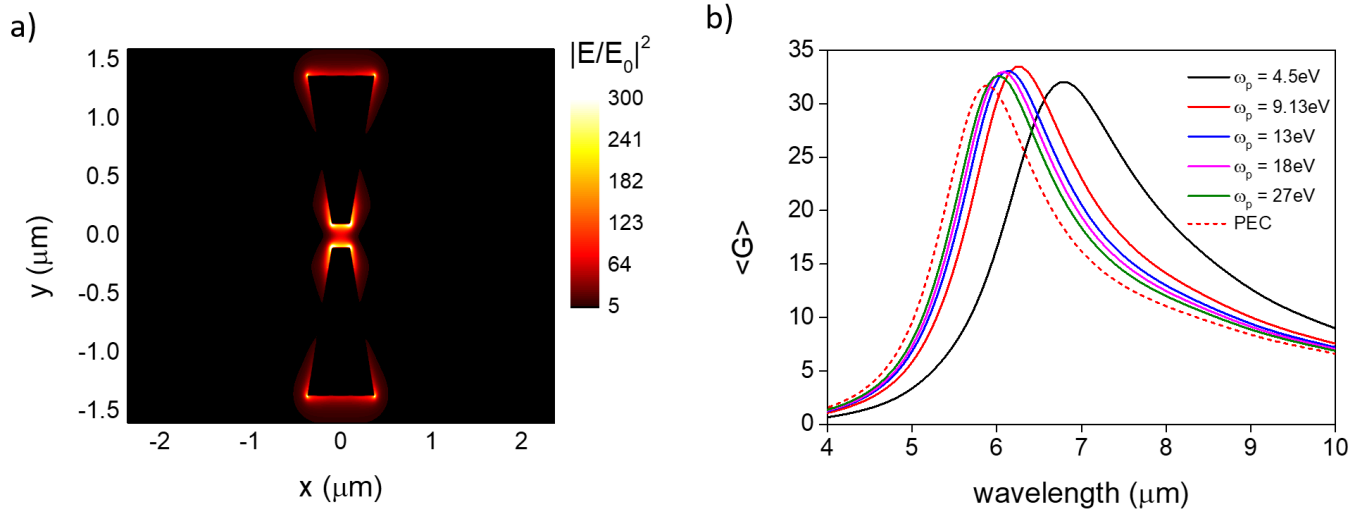
SUPPLEMENTARY NOTE 4: DEVICE RESISTANCE MODELING

The measured device resistance R_{D} as a function of symmetric gate voltage ($V_{\text{L}}=V_{\text{R}}$) is fitted considering $R_{\text{D}} = R_{\text{g}} + 2R_{\text{c}}$ and with fitting parameters the electron-hole mobilities, the residual local charge fluctuations n_0^* , the value of $\epsilon_{\text{hBN}}/d_{\text{hBN}}$, the Dirac voltage V_{D} (for charge neutrality) and the contact resistance R_{c} . The charge density distribution is obtained by the electrostatic calculations. We obtain $\mu_{\text{e}}(\mu_{\text{h}}) = 10,200$ (11,900) $\text{cm}^2\text{V}^{-1}\text{s}^{-1}$, $n_0^* = 1.72 \times 10^{11} \text{cm}^{-2}$, $V_{\text{D}} = 0.108$ V and $\epsilon_{\text{hBN}}/d_{\text{hBN}} = 0.222 \text{nm}^{-1}$. From the latter we extract hBN thickness $d_{\text{hBN}} = 15.7$ nm, in excellent agreement with the 15 nm used in this work, validating our approach. The contact resistance is described with a Gaussian distribution around the charge neutrality point¹⁵, with $R_{\text{c}} = 0.2$ (0.9) k Ω at high (low) doping (see Supplementary Figure 16).

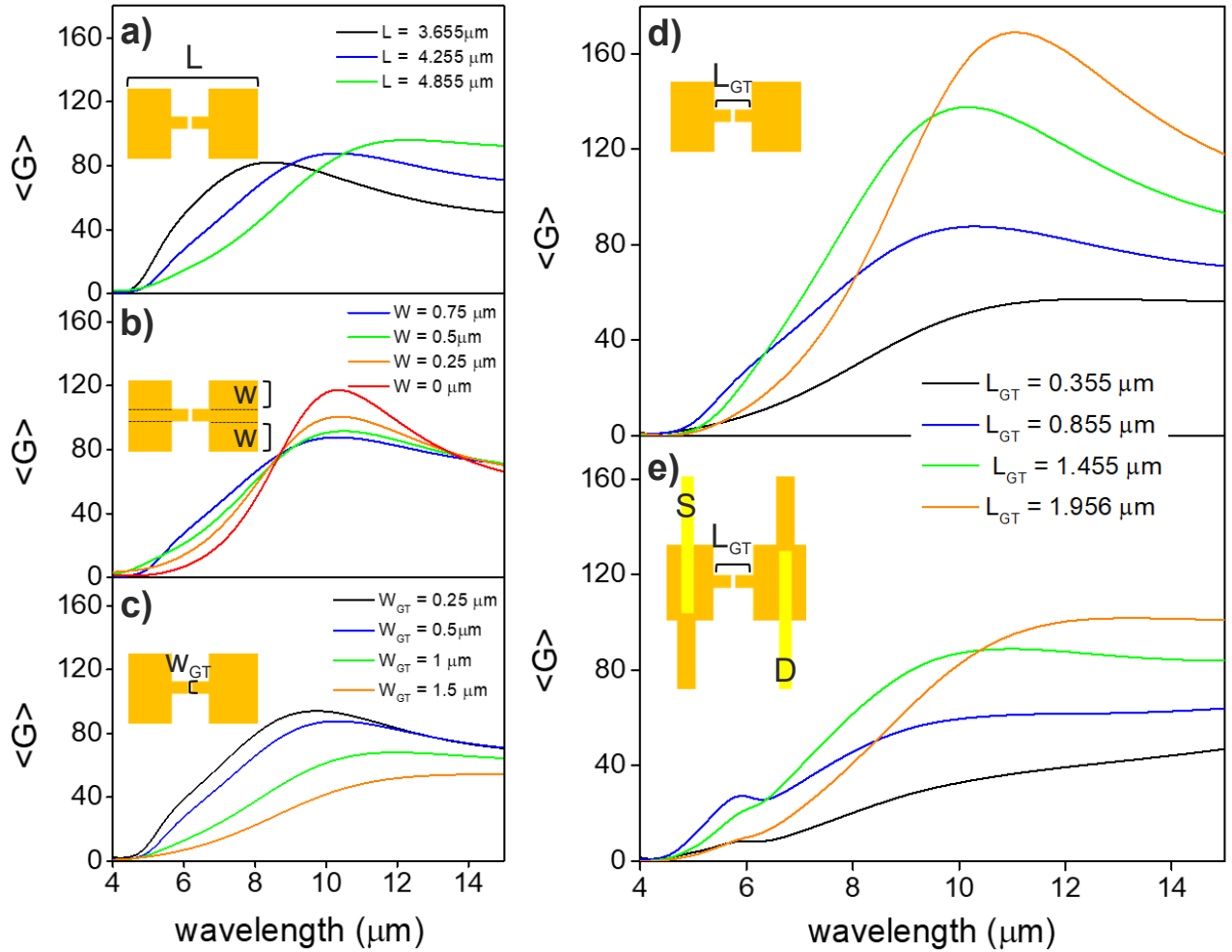
SUPPLEMENTARY NOTE 5: SPEED CALCULATIONS

The operation frequency of the photodetector is intrinsically related to the RC-time constant $\tau = R_{\text{D}}C_{\text{D}}$, with R_{D} as the total device resistance (see Supplementary Note 2 and 4) and C_{D} ¹⁶ as the total device capacitance given by $C_{\text{D}}^{-1} = C_{\text{q}}^{-1} + C_{\text{G}}^{-1}$, where C_{q} is the quantum graphene capacitance and C_{G} is the capacitance of the system as given by the normal formula for parallel-plate capacitors (for the case where $V_{\text{L}} = V_{\text{R}} = 0.5$ V and $\epsilon_{\text{hBN}}/d_{\text{hBN}} = 0.222 \text{nm}^{-1}$ as described in Supplementary Note 3 and 4). The operating speed is then described by the rate $f = (2\pi\tau)^{-1}$ and the rise time τ_{rise} , which is the time required for the photodetector to increase its output signal from 10% to 90% of the final steady-state output level. The rise time is calculated as $\tau_{\text{rise}} = \tau \times \ln(9) = (2\pi f)^{-1} \cdot \ln(9) = 0.35/f$.

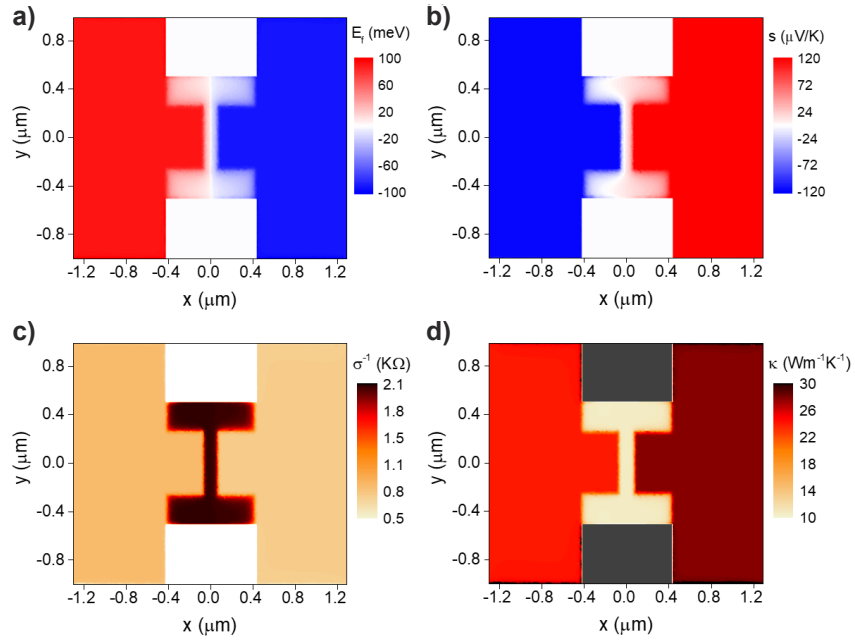
SUPPLEMENTARY FIGURES



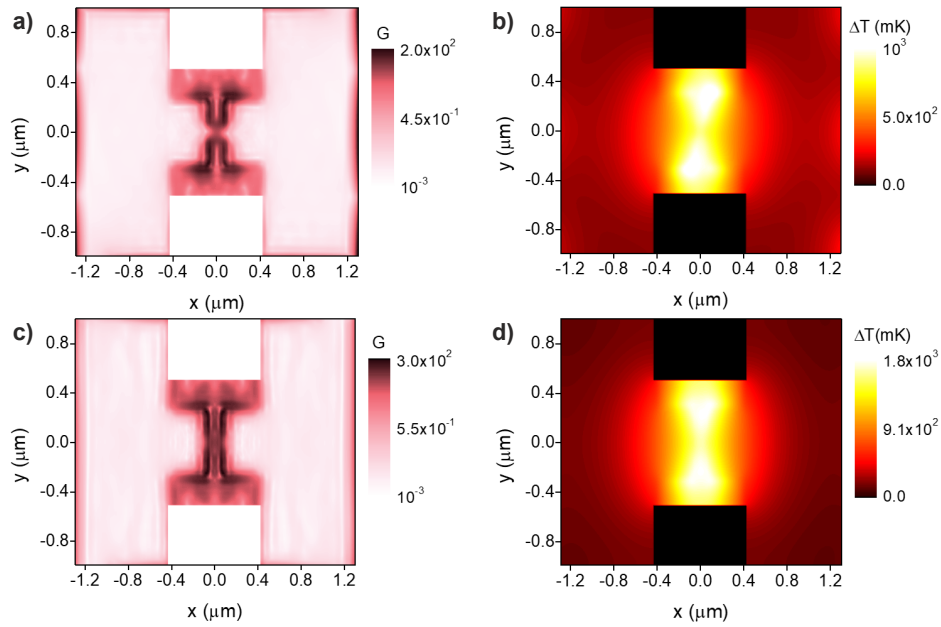
Supplementary Figure 2: Simulated spectral responsivities of **a)** TM and **b)** TE-polarization using different refractive indices configurations. Red solid line represents the case of wavelength-independent refractive indices for hBN ($n = 2.4$), SiO₂ ($n = 1.5$) and alumina ($n = 1.6$). For blue solid line we use the full dispersive optical model for the SiO₂ but not for the hBN, for which we use a wavelength independent refractive index of $n = 2.4$. Finally, for the black solid line we use full dispersive optical model for both SiO₂ and hBN.



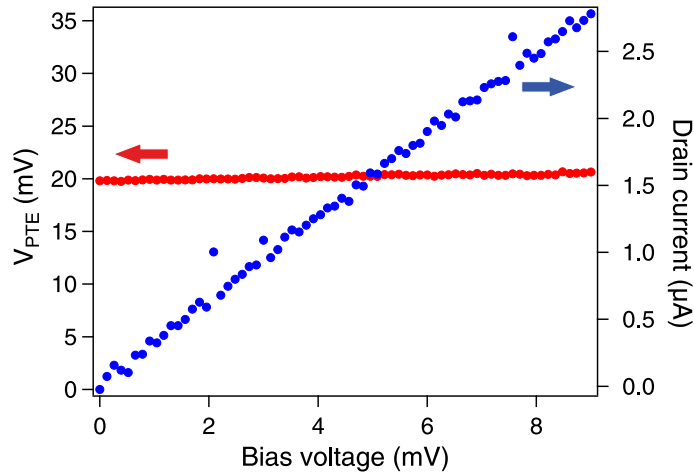
Supplementary Figure 3: Simulated $\langle G \rangle$, which is G averaged over the $0.2 \times 0.2 \mu\text{m}^2$ rectangular space around the center of the device ($x = 0, y = 0$, see Supplementary Figure 1a and Fig. 1a-b in the main text for axis definition), for TE-polarization as a function of the geometrical parameters of the H-shaped gates and additional metal components for the spectral range from 4 to 15 μm . In all cases, we use wavelength-independent refractive indices for hBN ($n = 2.4$), SiO_2 ($n = 1.5$) and alumina ($n = 1.6$), hence we can clearly observe the plasmonic response of the local gates without additional resonances coming from the phonon-polaritons of hBN or SiO_2 . In all cases the blue line corresponds to the response of the configuration where the varied parameter has the same value as the one it has in the experimental device. **a)** $\langle G \rangle$ values for different total lengths (L) of the local gates. We observe resonant peaks for all cases which shift to longer wavelengths as we increase L , a typical trend of metallic plasmonic resonators. In **b)** the width (W) of both extended parts of the H-shaped local gates is reduced starting from 0.75 μm (blue line) down to 0 μm (red line) while keeping L fixed at 4.255 μm . We clearly observe that the resonance peak is not shifted spectrally and also as W is reduced the spectral response converges to that of a dipole antenna (red line), thus proving the plasmonic behavior of the local gate geometry. In **c)**, we vary the gate tip width (W_{GT}) and notice that when increasing it, its plasmonic resonance is redshifted, its $\langle G \rangle$ amplitude drops down and the resonance becomes broader. In **d)** we vary the length of the gate tip (L_{GT}) while keeping L and W fixed at 4.255 μm and 0.75 μm respectively. We observe that the plasmonic response can be strongly tuned both spectrally and in amplitude when changing L_{GT} . Finally, in **e)** we add the extended electrodes to the local gates and also the source-drain contacts. We observe that these additional metal components further alter the response which is now lowered and broadened while we also notice an extra resonance peak at 5.8 μm . Note also that in the wavelength range of the hBN upper RB (6-7 μm) the optimum L_{GT} should be around 1.4 μm as shown in **d)** and **e)**. Due to the complexity of the H-shape configuration there are a large number of parameters that can strongly tune the plasmonic resonance of the local gates, meaning that we can further improve the device performance.



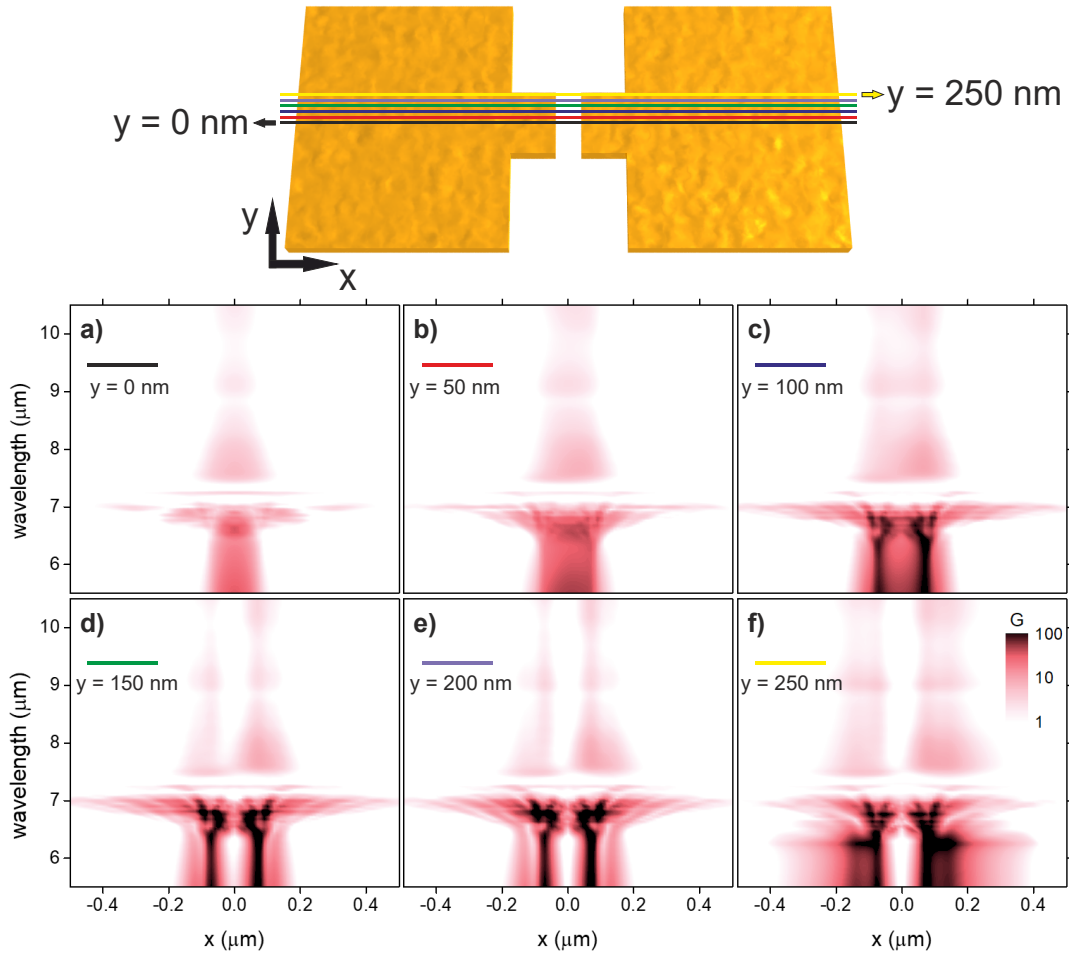
Supplementary Figure 4: Electrical and thermal properties of the graphene channel when applying 0.5 V (-0.5 V) to the left (right) gate. **a)** Calculated Fermi energy for the n -doped region (above left gate region) and p -doped region (above right gate region). In the gate gap, which is 155 nm, the Fermi energy drops to zero as we move from the n -doped to the p -doped gated region. Low values of the Fermi energy are also present in the graphene patches that extend 250 nm above and below the gate tips. **b)** Calculated Seebeck coefficient of the different graphene regions. As with the Fermi energy, the Seebeck coefficient drops to zero in the gate gap. Due to the difference in electron and hole mobilities, an imbalance is present in the weakly doped regions. Calculated **c)** resistivity and **d)** thermal conductivity. These two parameters are inversely proportional via the Wiedemann-Franz law.



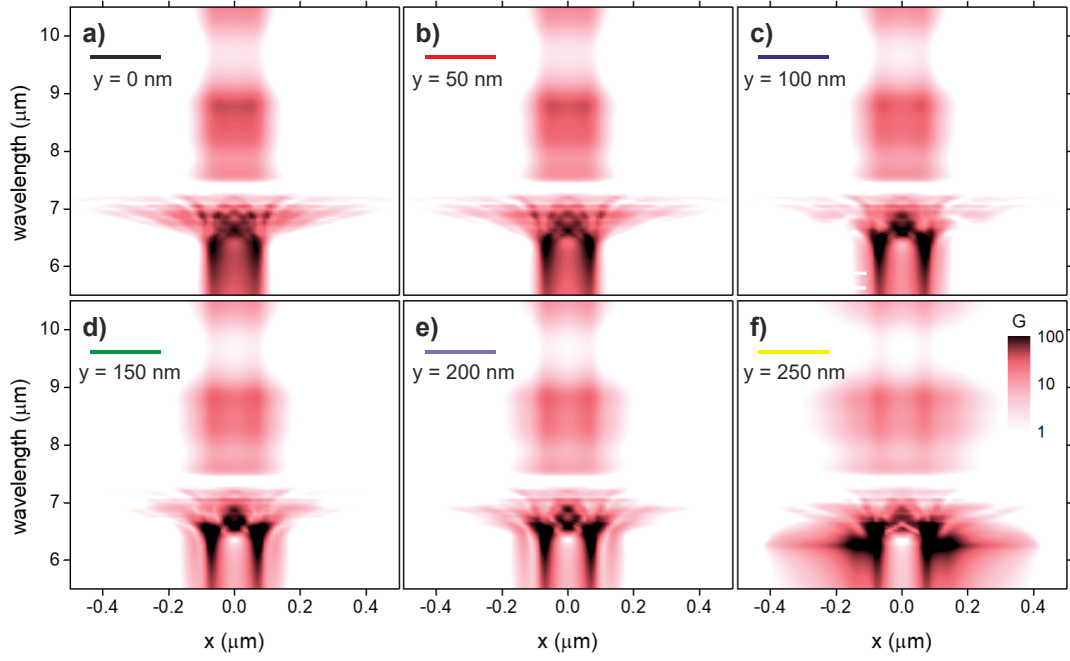
Supplementary Figure 5: a) G values in the graphene channel and b) temperature distribution for TM-polarization and c) and d) for TE-polarization respectively. For both cases the incident wavelength is $6.5 \mu\text{m}$. In both polarizations we observe the maximum values of G in an area of $0.4 \times 0.5 \mu\text{m}^2$ around the center of the pn -junction ($x=0, y=0$). The fact that in the rest of the graphene channel G values are 3-4 orders of magnitude lower compared to the ones above indicates efficient light focusing in the pn -junction. This enhanced absorption and low thermal conductivity (shown in Supplementary Figure 4d) across the pn -junction results in high temperature concentration, as is evident in b and d. Note that the difference in G values for the two cases corresponds to a difference in peak temperatures and finally in difference a responsivity as shown in Fig. 2 in the main text.



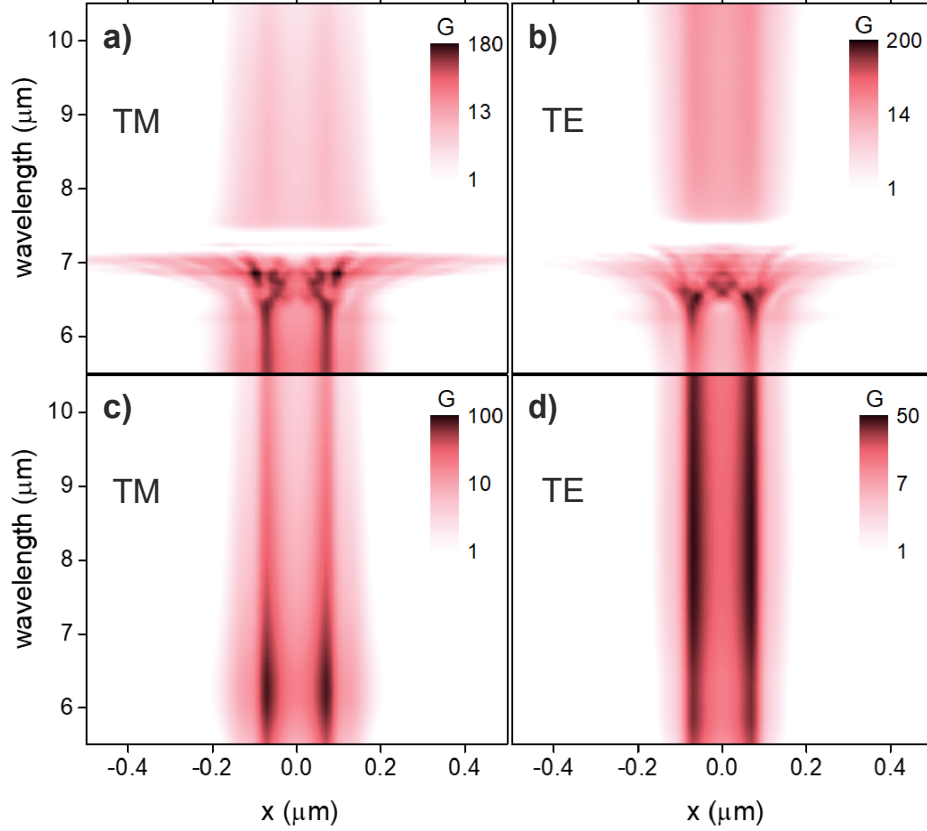
Supplementary Figure 6: Photovoltage (red) and source-drain current (blue) as a function of the graphene channel bias. We observe that the photocurrent remains constant, whereas the source-drain current increases linearly when increasing the bias voltage.



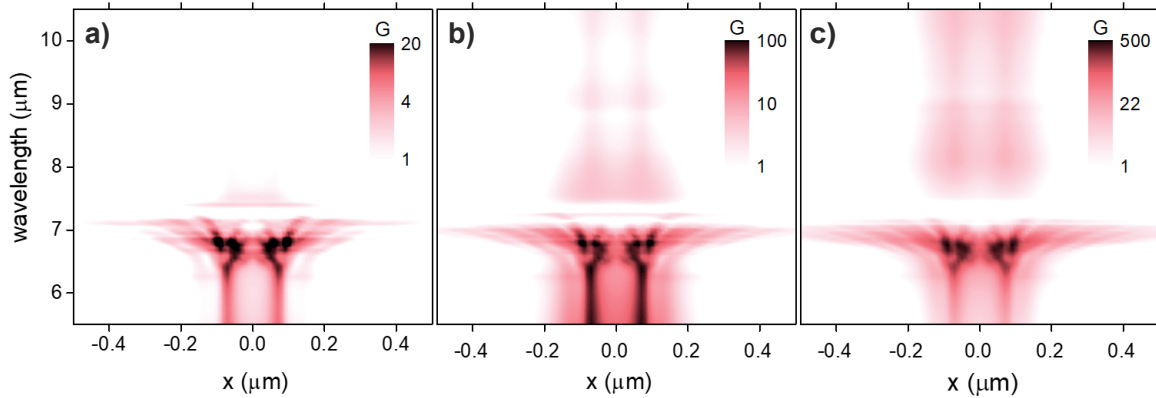
Supplementary Figure 7: G values for TM-polarization for different cross-sections across the gates along the source-drain direction (x direction, where $x = 0$ is located at the center of the gate gap. This gap is 155 nm. See Fig. 1a-b in the main text for axis definition) as a function of wavelength. **a)** corresponds to a linecut of G across the center of the gates ($y = 0$) as indicated in the top schematic. Then, the rest indicate linecuts of G **b)** 50 nm, **c)** 100 nm, **d)** 150 nm, **e)** 200 nm and **f)** 250 nm above the the center of the gates. The gates and antenna overlap at $y > 100$ nm (or $y < -100$ nm) away from the center ($y = 0$ nm). The upper antenna branch tip is located at $y = 100$ nm. We observe that the spatial pattern of G varies with y and that G increases significantly where these two metal regions overlap ($y > 100$ nm).



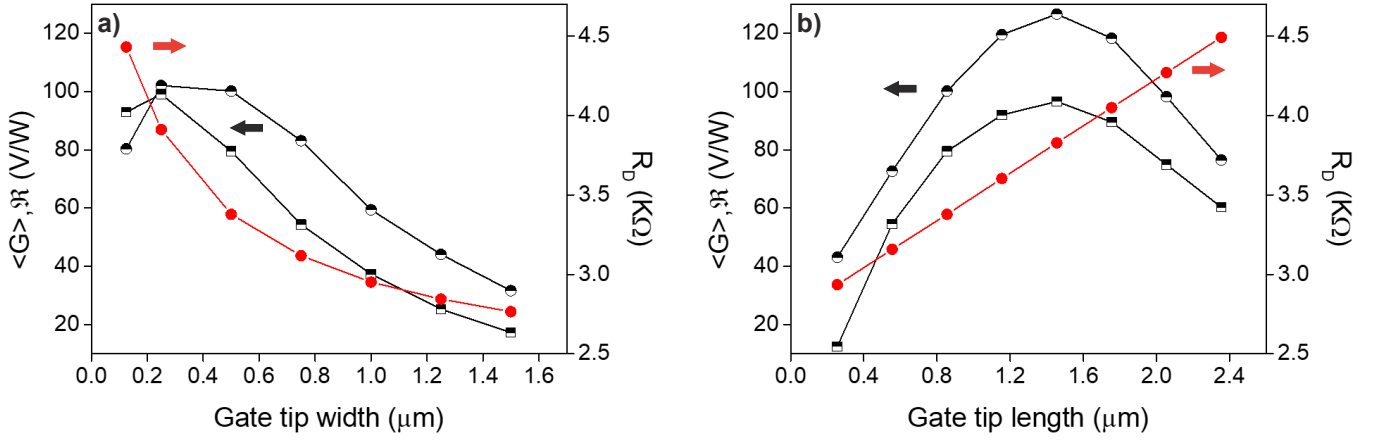
Supplementary Figure 8: Same as Supplementary Figure 7 but for TE-polarization. **a)** corresponds to a linecut of G across the center of the gates ($y = 0$). The rest indicate linecuts of G **b)** 50 nm, **c)** 100 nm, **d)** 150 nm, **e)** 200 nm and **f)** 250 nm above the the center of the gates. We observe that the spatial pattern of G varies while looking at it along different locations and that G greatly increases where the gate and antenna regions overlap ($y > 100$ nm).



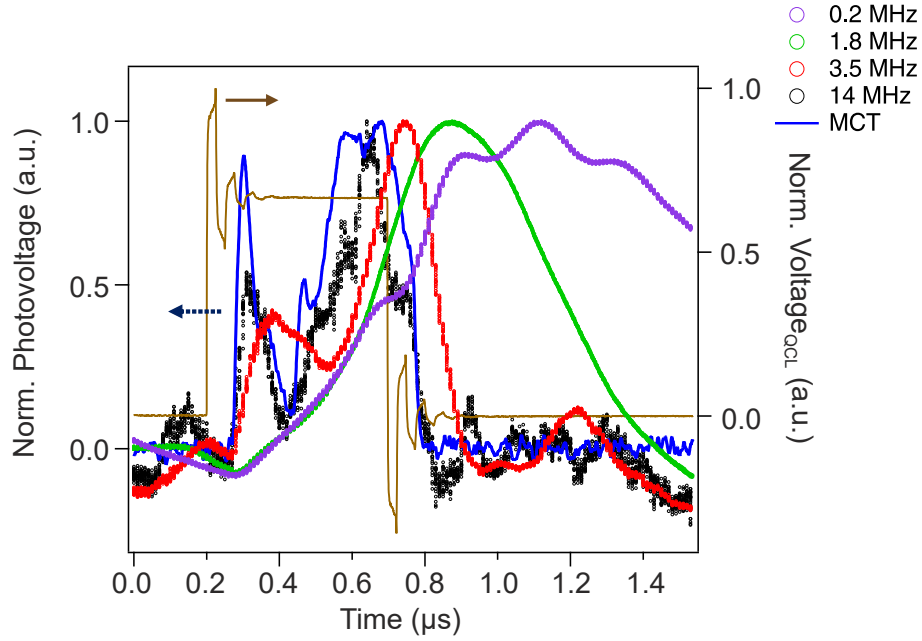
Supplementary Figure 9: Simulations of the absorption enhancement G along the source-drain direction (x direction and averaging over 500 nm in y direction, where $x = 0$ is located at the center of the gate gap. This gap is 155 nm. See Fig. 1a-b for axis definition) as a function of wavelength using different hBN and SiO_2 refractive indices configurations for TM (first column, **a**, **c**) and TE-polarization (second column, **b**, **d**). In **a**) and **b**), we use the full dispersive optical model for the hBN but not for the SiO_2 , for which we use a wavelength independent refractive index of $n = 1.5$. In **c**) and **d**), we use wavelength-independent refractive indices for hBN ($n = 2.4$), SiO_2 ($n = 1.5$) and alumina ($n = 1.6$).



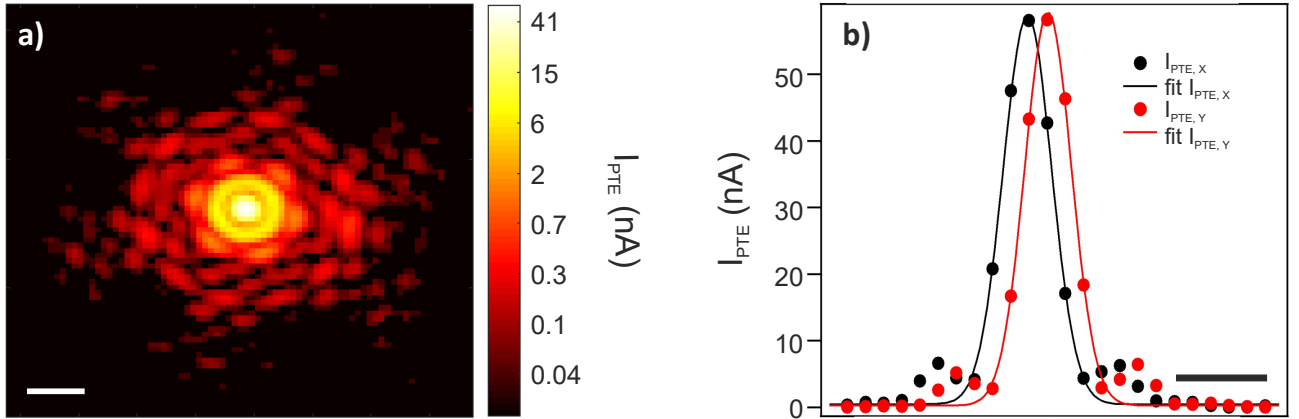
Supplementary Figure 10: 2D map of the simulated absorption enhancement G for TM-polarization along the source-drain direction (x direction and averaging over 500 nm in y direction, where $x = 0$ is located at the center of the gate gap. This gap is 155 nm. See Fig. 1a-b for axis definition) as a function of the wavelength (y -axis of the map) for total antenna lengths of **a**) $L = 1.8 \mu\text{m}$ (non-resonant antenna within hBN RB), **b**) $L = 2.7 \mu\text{m}$ (experimental/semi-resonant antenna within hBN RB) and **c**) $L = 4.8 \mu\text{m}$ (resonant antenna within hBN RB). For longer antennas, we notice that the LSPR of the antenna is strongly coupled with the hBN HPPs as observed in the G map in **c**.



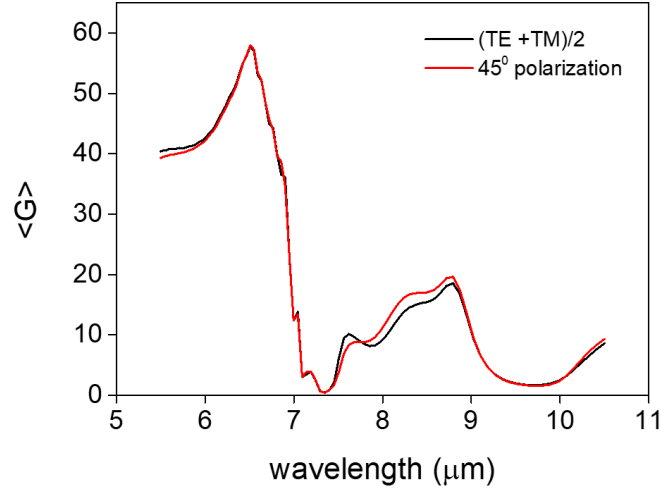
Supplementary Figure 11: Simulated responsivity (\mathfrak{R} in V/W, shown in black circles), $\langle G \rangle$ (defined in Supplementary Figure 3, black squares) and device resistance (R_D , shown in red circles) as a function of the geometric parameters of the H-shaped gates for TE-polarization at $\lambda = 6.5 \mu\text{m}$. **a)** Here we vary the tip width (W_{GT}). The gate tip length (L_{GT}) is fixed at 855 nm. We observe that W_{GT} correlates inversely with $\langle G \rangle$. We attribute this trend to the fact that the plasmonic response of the local gates decreases with the increase of W_{GT} as shown in Supplementary Figure 3c. The only case that does not follow this trend is the case of $W_{GT} = 0.125 \mu\text{m}$. In this case $\langle G \rangle$ decreases because W_{GT} is smaller than the $0.2 \mu\text{m}$ gap between the branches of the bow-tie antenna, thus preventing the gates and antenna from overlapping, while the overlap region precisely contains the highest values of G as shown in Supplementary Figure 8. As shown in Supplementary Figure 5, lower $\langle G \rangle$ values lead to smaller temperature gradients and consequently lower responsivities, since the responsivity is positively correlated with $\langle G \rangle$. Finally, as W_{GT} increases the width of the doped graphene channel also increases leading to smaller R_D . The interplay between responsivity and R_D gives the optimum W_{GT} case as shown in Fig. 4b in the main text. **b)** Here we vary L_{GT} . W_{GT} is fixed at 500 nm. We observe that $\langle G \rangle$ has a clear peak at $L_{GT} = 1.455 \mu\text{m}$. This trend is in excellent agreement with results in Supplementary Figure 3d-e, where the configuration using the same value of L_{GT} has the strongest response in the spectral position of the hBN upper RB. As in **a)**, the responsivity positively correlates with $\langle G \rangle$ and although R_D increases with L_{GT} , the optimum cases among all figures of merit (FOM) (Responsivity in V/W, A/W and NEP) remain roughly the same.



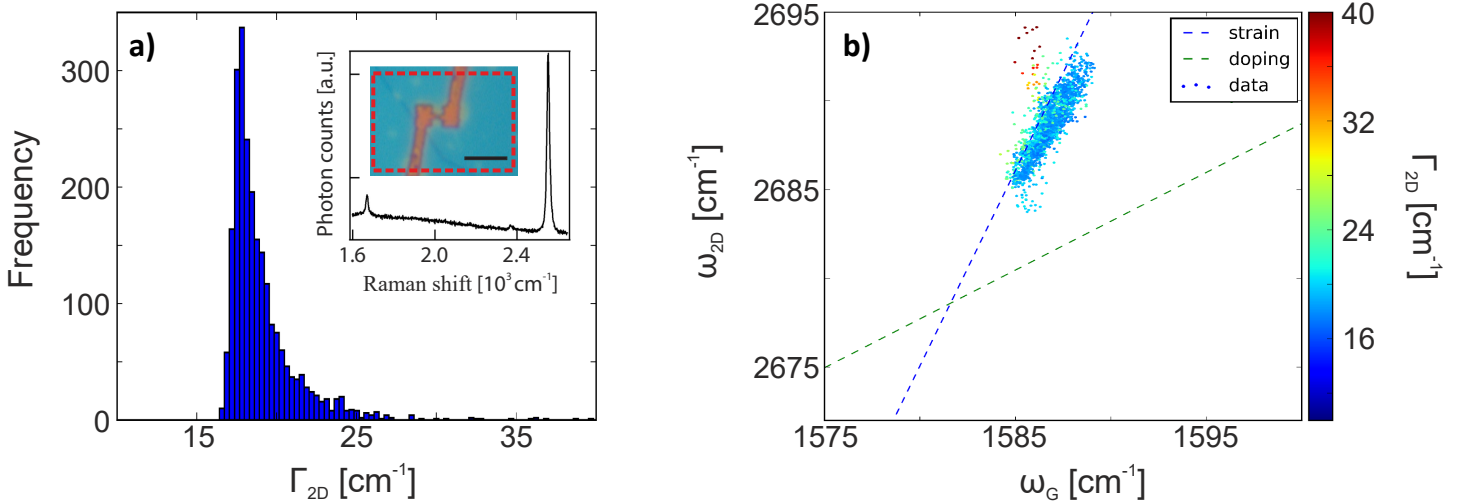
Supplementary Figure 12: Time-resolved photodetection traces for different operation bandwidths of the current amplifier. The speed increases (shorter rise time) with increasing bandwidth and the temporal response curve also resembles the photovoltage measured with the reference MCT detector as shown in blue solid line. The corresponding QCL voltage signal is shown in brown solid line.



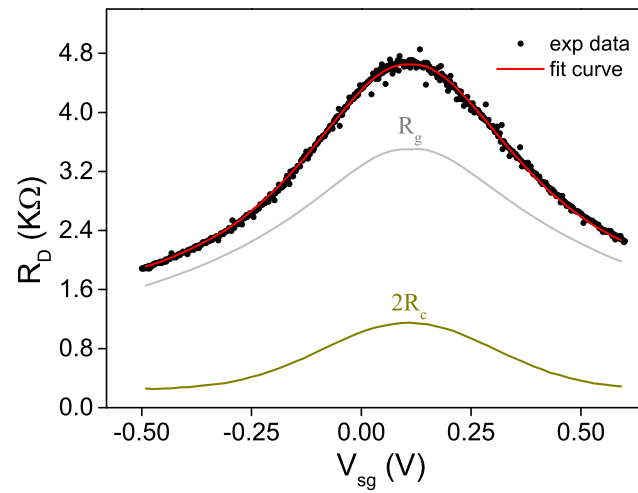
Supplementary Figure 13: Scanning photocurrent map at $\lambda = 6.6 \mu\text{m}$ with TE-polarization. **a)** Scanning photocurrent map (log scale) over the x and y scan directions. We observe an Airy beam pattern consisting of a central spot followed by several rings that contain a very small fraction of the total input power ($P_{\text{in}} = 13.7 \mu\text{W}$, $P_{\text{diff}} = 2.1 \mu\text{W}$). The white scale bar stands for $20 \mu\text{m}$. **b)** Linecuts of the map in **a)**, showing I_{PTE} across the x (black) and y (red) direction. The scale bar corresponds to $10 \mu\text{m}$. The dots represent the experimental I_{PTE} and the curves represent Gaussian fits. We obtain $w_{0,x} = 5.4 \mu\text{m}$ and $w_{0,y} = 5.2 \mu\text{m}$. The maximum responsivity value achieved was 27 mA/W (92 V/W) which corresponds to a NEP of $82 \text{ pW}/\sqrt{\text{Hz}}$, with a noise spectral density of $2.21 \text{ pA}/\sqrt{\text{Hz}}$ for a device resistance of $3.38 \text{ k}\Omega$ at the pn junction configuration given by $E_{\text{F,L}} = 85 \text{ meV}$, $E_{\text{F,R}} = -105 \text{ meV}$.



Supplementary Figure 14: Simulated absorption enhancement averaged over an area of $0.2 \times 0.2 \mu\text{m}^2$ around the device center ($x = y = 0$) for a 45 degree incident polarization respect to the bowtie antenna main axis (red curve) compared to the average TE and TM response (black curve). This result demonstrates that our system is linear and the photoresponse for an incident oblique polarization is the sum of the two components (TE and TM).



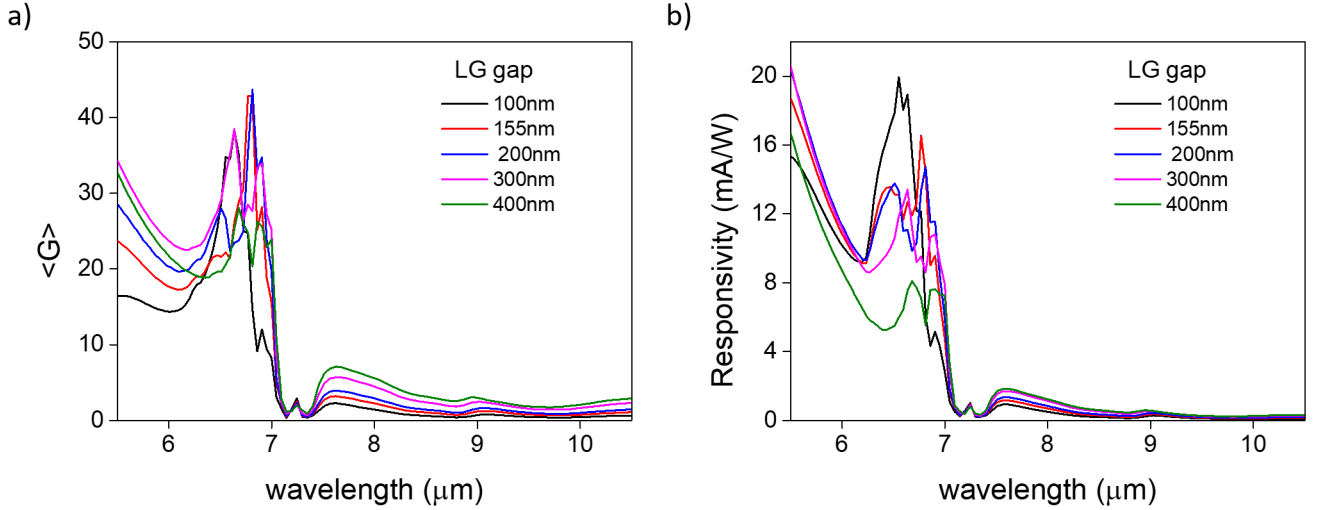
Supplementary Figure 15: Raman spectroscopy measurements of the 2D stack on top of the metal gates. **a)** Histogram plot of the full-width-half-maximum of the graphene 2D peak (Γ_{2D}) across a region of $\sim 13.5 \times 9 \mu\text{m}^2$ as shown in the inset in red dashed line. The inset scale bar corresponds to $4.2 \mu\text{m}$. The mean Γ_{2D} is $\sim 18 \text{ cm}^{-1}$, which reflects the high quality of the monolayer graphene encapsulated in hBN. The inset plot corresponds to the usual spectrum obtained in these measurements for single layer graphene. **b)** 2D frequency peak (ω_{2D}) as a function of G frequency peak (ω_G) obtained from the Raman map in panel **a**, where the colorbar corresponds to the Γ_{2D} of the measured map. We obtained low doping values consistent with transport measurements (see Supplementary Figure 16) and modest strain values¹⁷.



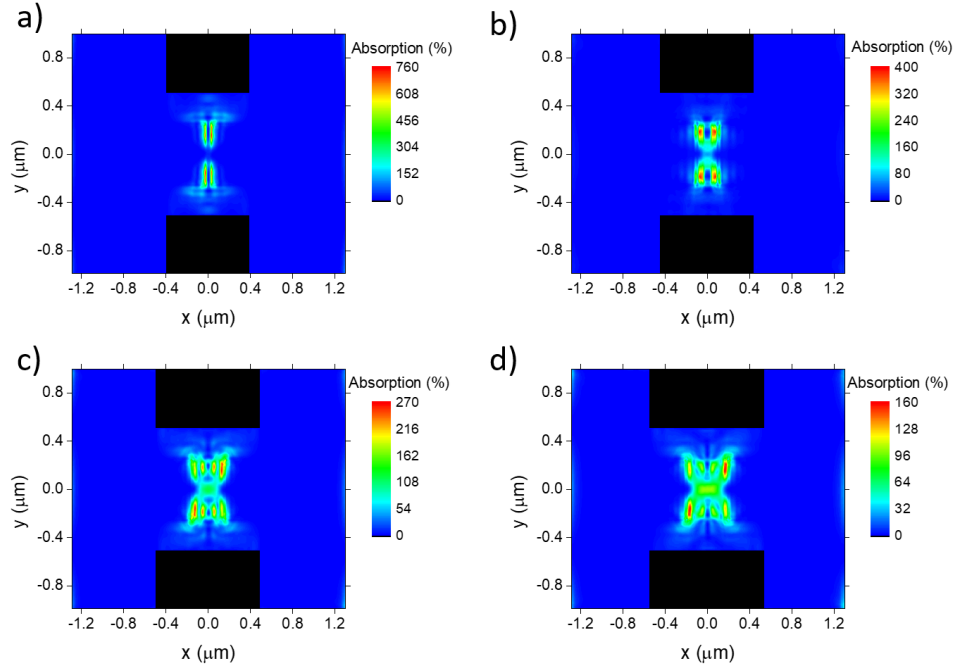
Supplementary Figure 16: Measured device resistance (R_D) as a function of the two gate voltages (V_L and V_R) both sweeping at the same voltage (V_{sg}). We fit the resistance curve using the model described in Supplementary Note 4. Contributions of both contact resistance ($2R_c$) and graphene channel resistance (R_g) to R_D for all V_{sg} are also presented.

SUPPLEMENTARY DISCUSSION 1: GATE GAP EFFECT ON THE PHOTORESPONSE

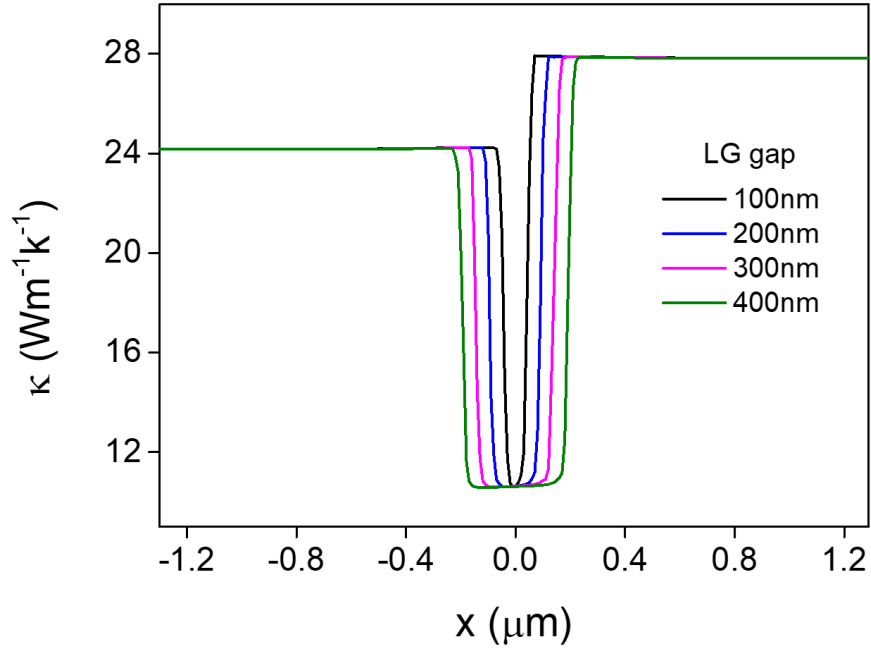
We perform a set of simulations for the TM-polarization case changing the gate gap between 100 nm and 400 nm to observe its effect on the photoresponse performance. Specifically, as shown in Supplementary Figure 17a, the averaged spectral absorption enhancement profile $\langle G \rangle$ is different for all cases, especially for the case of 400 nm gap where the peak values are much lower than the rest of the cases. This indicates that the interference patterns of hBN HPPs are affected by the gate gap. In terms of responsivity (Supplementary Figure 17b), these differences become more pronounced, as we can observe a clear optimum case for the 100 nm gap (the lower limit in terms of fabrication), while the performance of the 400 nm gap case is by far inferior. To elucidate this behaviour, we examine the spatial profiles of the parameters that affect the device performance, such as absorption (Supplementary Figure 18), hot electron temperature distribution (Supplementary Figure 20) and PTE voltage density (Supplementary Figure 22) for the cases of 100, 200, 300 and 400 nm gap at the wavelength where the maximum responsivity is observed for each case. Absorption profiles in Supplementary Figure 18 clearly explains the trends in Supplementary Figure 17a. As the overlap between the bowtie antenna and the tips of the gates decreases (by increasing the gap of the latter), both relative strength and spatial pattern is strongly altered. For 100 nm gap we observe large values of absorption, which however are concentrated beyond the gate edges where thermal conductivity is high (see Supplementary Figure 19). For 400 nm gap, on the other hand (no gate-antenna overlap) the absorption is strongly diminished but also more uniformly distributed inside the gap where thermal conductivity is low (see Supplementary Figure 19). These two conflicting attributes determine the temperature distribution as presented in Supplementary Figure 20. Indeed, as we increase the gap the differences we observe in peak temperature are smaller than those in absorption. Nonetheless, if we consider that for larger gaps the temperature peak is well centred within the gap area where the Seebeck coefficient (Supplementary Figure 21) has low values, we understand why PTE voltage generation diminishes for larger gaps, as is evident in Supplementary Figure 22. Taking also into account that resistance increases for larger gaps (see Supplementary Figure 23) explains the performance of our device as presented in Supplementary Figure 17b.



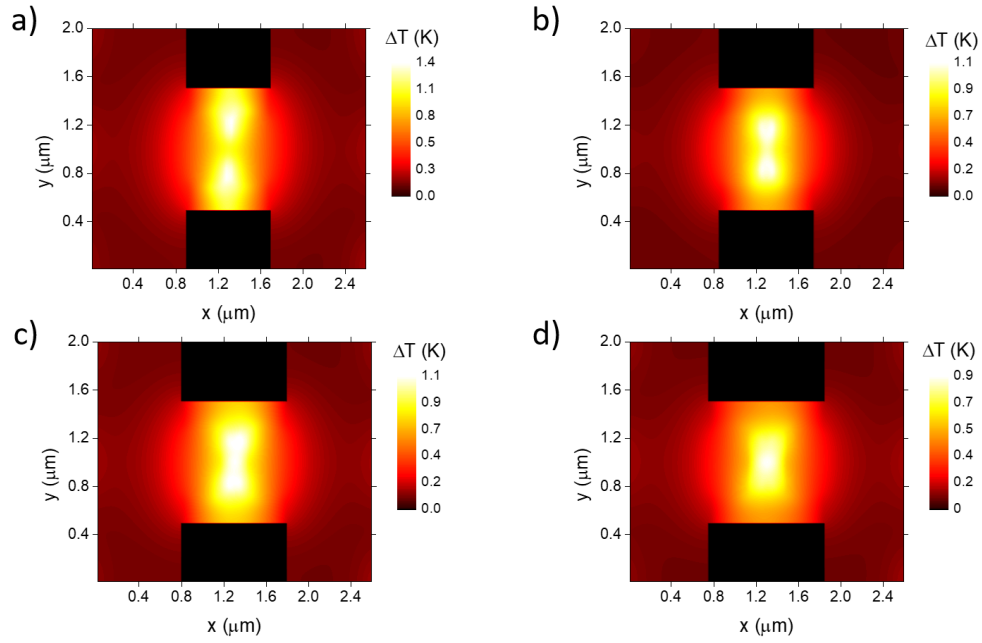
Supplementary Figure 17: **a)** Simulated $\langle G \rangle$ (averaged over an area $0.5 \times 0.5 \mu\text{m}^2$ around the device center $x = y = 0$) for devices with different gate gaps with incident TM-polarization, **b)** corresponding spectral responsivity. The gate voltages configuration is $V_{L,R} = \pm 0.5 \text{ V}$ ($E_F = \pm 100 \text{ meV}$).



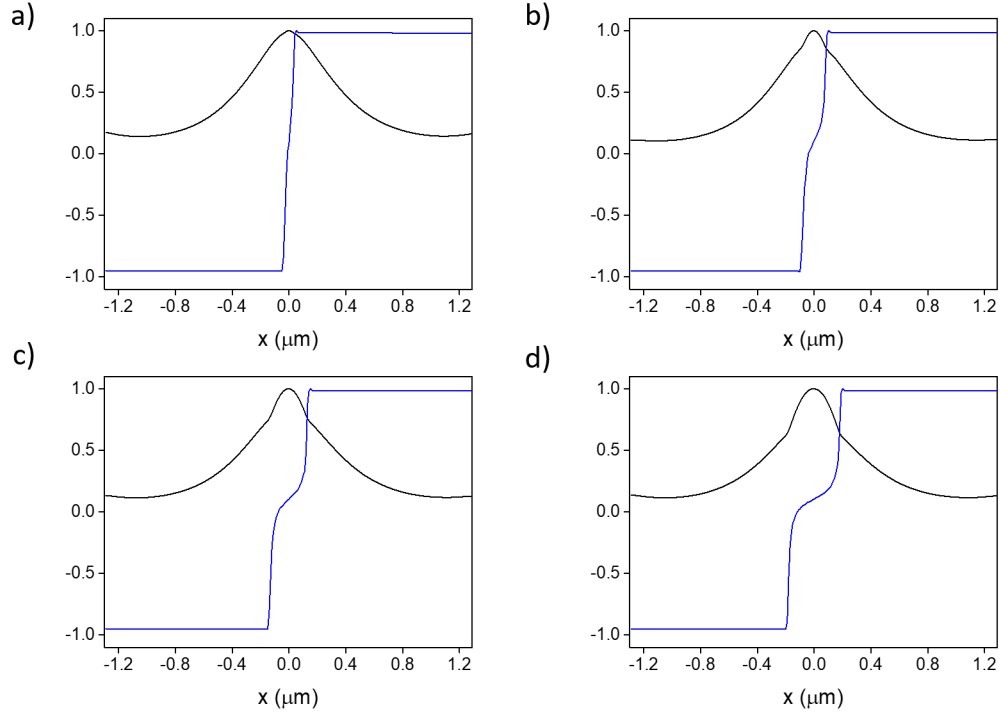
Supplementary Figure 18: Absorption distribution in graphene for TM-polarization with gate gap of **a)** 100 nm, **b)** 200 nm, **c)** 300 nm and **d)** 400 nm. The peak wavelength is chosen in each case ($\lambda = 6.551, 6.767, 6.636$ and $6.679 \mu\text{m}$ respectively). The gate voltage configuration is the same as displayed in Supplementary Figure 17.



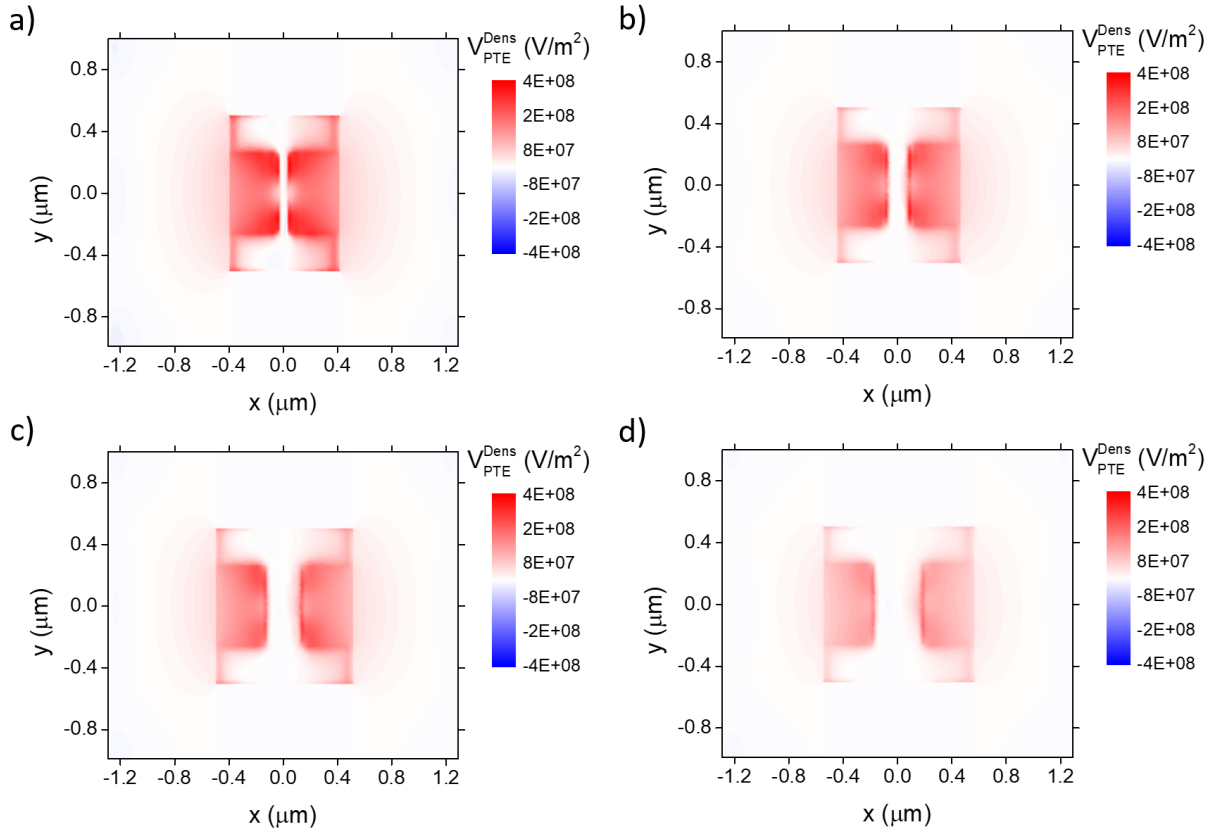
Supplementary Figure 19: Thermal conductivity across the graphene channel (at $y = 0$) for the devices with the different gate gaps and gate voltage configuration $V_{L,R} = \pm 0.5 \text{ V}$ ($E_F = \pm 100 \text{ meV}$). The small imbalance is due to the difference in electron and hole mobilities.



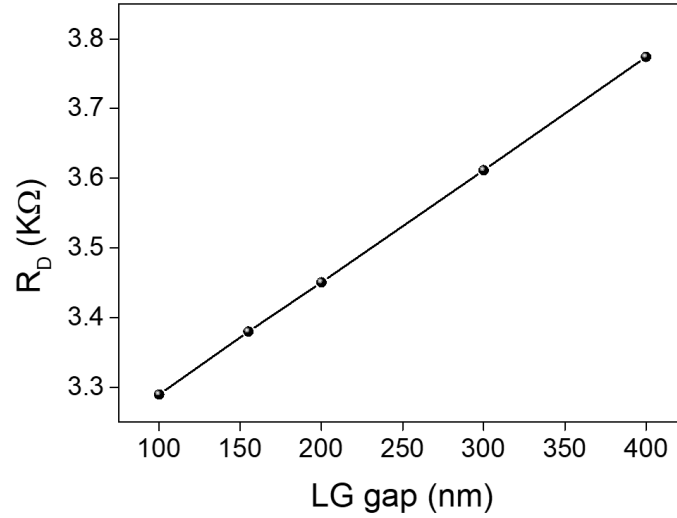
Supplementary Figure 20: Temperature distribution in graphene electrons for TM-polarization with gate gap of **a)** 100 nm, **b)** 200 nm, **c)** 300 nm and **d)** 400 nm. The peak wavelength chosen in each case is explained in Supplementary Figure 18. The gate voltages configuration is $V_{L,R} = \pm 0.5$ V ($E_F = \pm 100$ meV).



Supplementary Figure 21: Normalized ΔT (black lines) and Seebeck coefficient (blue lines) across the graphene channel (at $y = 0$) for TM-polarization with gate gap of **a)** 100 nm, **b)** 200 nm, **c)** 300 nm and **d)** 400 nm. The peak wavelength chosen in each case is explained in Supplementary Figure 18. The gate voltages configuration is $V_{L,R} = \pm 0.5$ V ($E_F = \pm 100$ meV).



Supplementary Figure 22: PTE voltage density distribution in graphene for TM-polarization with gate gap of a) 100 nm, b) 200 nm, c) 300 nm and d) 400 nm. The peak wavelength chosen in each case is explained in Supplementary Figure 18. The gate voltages configuration is $V_{\text{L,R}} = \pm 0.5$ V ($E_{\text{F}} = \pm 100$ meV).



Supplementary Figure 23: Device resistance as function of the gate gap for the gate configuration $V_{\text{L,R}} = \pm 0.5$ V ($E_{\text{F}} = \pm 100$ meV).

SUPPLEMENTARY REFERENCES

- ¹ Woessner, A. et al. [Highly confined low-loss plasmons in graphene–boron nitride heterostructures](#). *Nature Mater.* **14**, 421–425 (2015).
- ² Kischkat, J. et al. [Mid-infrared optical properties of thin films of aluminum oxide, titanium dioxide, silicon dioxide, aluminum nitride, and silicon nitride](#). *Applied Optics* **51**, 6789–6798 (2012).
- ³ Hanson, G. W. [Quasi-transverse electromagnetic modes supported by a graphene parallel-plate waveguide](#). *Journal of Applied Physics* **104** (2008).
- ⁴ Johnson, P. & Christy, R. [Optical constants of the noble metals](#). *Phys. Rev. B* **6**, 4370–4379 (1972).
- ⁵ Song, J. C. W., Rudner, M. S., Marcus, C. M. & Levitov, L. S. [Hot carrier transport and photocurrent response in graphene](#). *Nano Letters* **11**, 4688–4692 (2011).
- ⁶ Castilla, S. et al. [Fast and Sensitive Terahertz Detection Using an Antenna-Integrated Graphene pn Junction](#). *Nano Letters* **19**, 2765–2773 (2019).
- ⁷ Tielrooij, K.-J. et al. [Out-of-plane heat transfer in van der Waals stacks through electron–hyperbolic phonon coupling](#). *Nature Nanotechnology* **13**, 41–46 (2018).
- ⁸ Tan, Y. W. et al. [Measurement of Scattering Rate and Minimum Conductivity in Graphene](#). *Phys. Rev. Lett.* **99**, 246803 (2007).
- ⁹ Zuev, Y., Chang, W. & Kim, P. [Thermoelectric and Magnetothermoelectric Transport Measurements of Graphene](#). *Phys. Rev. Lett.* **102**, 96807 (2009).
- ¹⁰ Cutler, M. & Mott, N. F. [Observation of anderson localization in an electron gas](#). *Physical Review* **181**, 1336–1340 (1969).
- ¹¹ Soavi, G. et al. [Broadband, electrically tunable third-harmonic generation in graphene](#). *Nature Nanotechnology* **13**, 583–588 (2018).
- ¹² Kittel, C. [Introduction to Solid State Physics](#) (Wiley, 2004), 8 edn.
- ¹³ McPherson, J. W., Kim, J., Shanware, A., Mogul, H. & Rodriguez, J. [Trends in the ultimate breakdown strength of high dielectric-constant materials](#). *IEEE Transactions on Electron Devices* **50**, 1771–1778 (2003).
- ¹⁴ Kim, K. K. et al. [Synthesis and characterization of hexagonal boron nitride film as a dielectric layer for graphene devices](#). *ACS Nano* **6**, 8583–8590 (2012).
- ¹⁵ Xia, F., Perebeinos, V., Lin, Y. M., Wu, Y. & Avouris, P. [The origins and limits of metal-graphene junction resistance](#). *Nature Nanotechnology* **6**, 179–184 (2011).
- ¹⁶ Mišković, Z. L. & Upadhyaya, N. [Modeling electrolytically top-gated graphene](#). *Nanoscale Research Letters* **5**, 505–511 (2010).
- ¹⁷ Lee, J. E., Ahn, G., Shim, J., Lee, Y. S. & Ryu, S. [Optical separation of mechanical strain from charge doping in graphene](#). *Nature Communications* **3**, 1024–1028 (2012).